

Dissertation
submitted to the
Combined Faculties for the Natural Sciences and for Mathematics
of the Ruperto-Carola University of Heidelberg, Germany
for the degree of
Doctor of Natural Sciences

presented by
Sasithorn Chotewutmontri (M.Sc. Biotechnology)
born in Nakhon Phanom, Thailand
Oral-examination:

Genome integration structures and genotype variants of oncogenic human papillomavirus types HPV16 and HPV68 in cervical carcinoma-derived cell lines, cervical precursor lesions and carcinomas

Referees:

Prof. Dr. Elisabeth Schwarz
Prof. Dr. Gabriele Petersen

First of all, I would like to express my utmost gratitude to my mentor, Prof. Dr. Elisabeth Schwarz, for giving me the opportunity to work in this interesting field, for her constant support and valuable advice in both scientific and personal matters, for her critical guidance, endless patience and the wealth of suggestions and comments in the completion of this dissertation, and for translating the Abstract into the German Zusammenfassung. At the same time, I would like to thank Prof. Dr. Gabriele Petersen for kindly being my second supervisor, and PD Dr. Stefan Wiemann for being a member of my advisory committee, and both of them for their attendance and discussions in my progress reports. Furthermore, I would like to thank PD Dr. Anne Régnier-Vigouroux and Prof. Dr. Rainer Zawatzky for kindly being my third and fourth examiners.

I would also like to thank our partners in the Cancéropôle du Grand-Est - DKFZ collaboration project, especially, Dr. Véronique Dalstein and Prof. Dr. Christine Clavel (Université de Reims Champagne-Ardenne, Reims, France), and Dr. Maëlle Saunier and Dr. Jean-Luc Prétet (Université de Franche-Comté, Besançon, France) for providing the clinical samples for this work.

I would like to specially thank Dr. Bo Xu for giving me the opportunity to collaborate in the development of the ASP16 strategy and for his advice and support. I also would like to give my thanks to Ursula Klos, Monika Frank-Stöhr, Ilona Braspenning-Wesch and Birgit Hub for laboratory advice and assistance, and for their friendship. Together, they created the most pleasant and friendly working environment I have ever known. Furthermore, I would like to thank Prof. Dr. Frank Rösl and the members of the Department of Viral Transformation Mechanisms at DKFZ for their support and friendliness.

I would like to thank Andreas Hunziker for Sanger sequencing, and Dr. Stephan Wolf for Roche/454 GS-FLX pyrosequencing and related discussions.

I would like to thank Dr. Kajohn Boonrod, Ria Kretzer and Jochen Kretzer for their constant support and encouragement in difficult times, and for being my second family.

Finally, I would like to dedicate my thesis to my beloved parents, sister Joy and brother Non. Without them, I would not have achieved this far.

Für die Entstehung von Gebärmutterhalskrebs (Zervixkarzinom) ist eine persistierende Infektion mit humanen Papillomviren (HPV) der Hochrisiko-Gruppe die entscheidende Voraussetzung. Häufig erfolgt später eine Integration der viralen DNA in das Genom der Wirtszellen. Durch die Integration wird meistens das virale E2-Gen zerstört oder deletiert, wodurch es zur Deregulation der in der „upstream regulatory region“ (URR) gestarteten Transkription der viralen Onkogene E6 und E7 kommt. Einen zusätzlichen Einfluss auf den Mehrstufenprozess der Zervixkarzinogenese kann die integrierte HPV-DNA ausüben, indem wichtige zelluläre Gene durch Insertionsmutagenese verändert werden. HPV16 ist der häufigste und HPV68 ein seltener Hochrisiko-HPV-Typ. HPV16 kommt in ungefähr 55 % der weltweit untersuchten Zervixkarzinome vor, HPV68 in weniger als 1 %. In dieser Arbeit wurden in Zervixkarzinom-Zelllinien und in klinischen Proben von Zervixabstrichen die DNA-Strukturen von HPV68 bestimmt sowie HPV16-Integration und Sequenzen des E1-E2-Genbereichs mittels der neuen ASP16-Strategie („amplification selection pyrosequencing of HPV16“) analysiert.

HPV68 wird in zwei Subtypen, a und b, unterteilt. Besonderes Merkmal von HPV68b ist das Vorkommen als integrierte DNA in der Zervixkarzinom-Zelllinie ME180. Die davon abgeleitete Linie ME180R, resistent gegenüber Wachstumsinhibition durch Tumor-Nekrose-Faktor alpha (TNFalpha), weist partielle Deletionen in der integrierten HPV68b-DNA auf. In dieser Arbeit wurden die kompletten Strukturen der integrierten HPV68b-DNA in ME180 und ME180R bestimmt. ME180-Zellen enthalten zwei unvollständige Kopien von HPV68b, die in einer einzigartigen Kopf-Kopf-Anordnung integriert vorliegen. Durch Selektion von zwei neuen TNFalpha-resistenten ME180-Sublinien konnte gezeigt werden, dass die Umlagerungen und partiellen Deletionen der integrierten HPV68b-DNA in ME180R nicht wesentlich für den TNFalpha-resistenten Phänotyp sind. Aus einer CIN2-Krebsvorstufe (zervikale intraepitheliale Neoplasie Grad 2) wurden ein komplettes und ein mutiertes HPV68b-Genom isoliert, kloniert und sequenziert. Das mutierte Genom, das eine 1,2 kb große Deletion im E1-Gen aufweist, liegt wahrscheinlich integriert vor. In elf HPV68-positiven klinischen Proben wurde das virale Genom durch partielle URR-Sequenzierung untersucht. Eine Probe enthielt HPV68a, die anderen zehn Proben enthielten HPV68b-Varianten, von denen neun vorher unbekannt waren. Die Ergebnisse zeigen, dass HPV68b häufiger als HPV68a und in vielen molekularen Varianten vorkommt.

Die ASP16-Strategie wurde entwickelt zur gleichzeitigen Bestimmung von HPV16-Integrationsstellen in einer Vielzahl von klinischen DNA-Proben. ASP16 besteht aus vier Hauptschritten: GenomePlex Gesamtgenom-Amplifikation, Anreicherung von HPV16 E1-E2-Sequenzen, Roche/454 GS-FLX Pyrosequenzierung, und Daten-Analyse. In dieser Arbeit wurden Computerprogramme zur ASP16-Datenanalyse entwickelt und eingesetzt. Die ASP16-Strategie wurde weiter optimiert und zur Analyse von 25 HPV16-positiven Proben eingesetzt. Es wurden längere Einzelsequenzen als vorher sowie eine Sequenzabdeckung von 89 % erreicht. HPV16-Integrationsstellen konnten in 3 von 4 Zelllinien und in 6 von 21 klinischen Proben identifiziert werden. Die in den klinischen Proben identifizierten HPV16-Integrationsstellen liegen alle in oder in der Nähe von zellulären Proto-Onkogenen oder Tumorsuppressorgenen. Diese Befunde unterstützen die Vermutung, dass die HPV-Integration durch Veränderungen krebsrelevanter zellulärer Gene zur Zervixkarzinogenese beitragen kann. Die hohe Sequenzabdeckung im E1-E2-Bereich ermöglichte auch die Bestimmung von HPV16-Varianten. Die ASP16-Strategie ist die erste Methode, die „next generation sequencing“-Technologien mit der Bestimmung von HPV-Integrationsstellen im Multiplex-Format kombiniert. ASP16 ermöglicht damit die serienmäßige Analyse von HPV16-Integrationsstellen in klinischen Proben und liefert gleichzeitig E1-E2-Sequenzen zur Bestimmung von Mutationen und Varianten.

Persistent infection with high-risk human papillomavirus (hr-HPV) is essential for cervical carcinogenesis, and is frequently followed by integration of the viral DNA into the host genome. Upon integration, the viral E2 gene is usually disrupted or deleted leading to deregulated transcription of the E6/E7 oncogenes from the upstream regulatory region (URR). Integrated HPV DNA may also affect critical cellular genes through insertional mutagenesis, which can contribute to the multi-step process of cervical carcinogenesis. HPV16 is the most frequent and HPV68 is a rare hr-HPV type, present in about 55% and less than 1% of cervical carcinomas worldwide, respectively. In this work, HPV68 DNA structures in cervical carcinoma cell lines and clinical samples were analyzed. HPV16 integration and E1-E2 sequences were studied using the novel “amplification selection pyrosequencing of HPV16” (ASP16) strategy.

HPV68 is divided into two subtypes, a and b. A hallmark of HPV68b is its presence as integrated DNA in the cervical carcinoma cell line ME180. In the mutant cell line ME180R, selected for resistance to growth inhibition by tumor-necrosis-factor alpha (TNFalpha), partial deletions in the integrated HPV68b DNA had been detected. In this study, the complete structures of the integrated HPV68b in ME180 and ME180R have been determined. ME180 cells contain two disrupted HPV68b copies, integrated in a unique head-to-head arrangement into chromosome 18q21. By selection of new TNFalpha-resistant ME180 sub-lines, it was found that the rearrangements and partial deletions of HPV68b in ME180R are unnecessary for the TNFalpha-resistance phenotype. In addition, a full-length and a mutant HPV68b genome were isolated from a cervical intraepithelial neoplasia grade 2 (CIN2) precursor lesion, cloned and completely sequenced. The mutant genome carrying a 1.2-kb deletion in the E1 gene is probably integrated. Based on partial URR sequences, ten HPV68b variants, nine of them new, and one HPV68a variant have been identified in eleven clinical samples, suggesting that HPV68b is more widely distributed than HPV68a and is present in a multitude of molecular variants.

ASP16 was developed for simultaneous determination of HPV16 integration junctions in multiple clinical DNA samples. It consists of four main steps: GenomePlex whole genome amplification, HPV16 E1-E2 sequence enrichment, Roche/454 GS-FLX pyrosequencing, and data analysis. In this work, computer programs for ASP16 data analysis were developed and applied. The ASP16 strategy was further optimized and used for the analysis of 25 HPV16-positive samples. The optimized ASP16 delivered longer sequence read lengths and 89% average sequence coverage. HPV16 integration junctions were identified in 3 out of 4 cell lines, and 6 out of 21 clinical samples. The HPV16 integration sites identified in the clinical samples are all located near cellular proto-oncogenes or tumor suppressor genes, supporting the assumption that HPV integration contributes to cervical carcinogenesis by altering cancer-relevant cellular genes. The high E1-E2 sequence coverage also allowed HPV16 variant assignments. Altogether, the ASP16 strategy, which is the first method combining next generation sequencing technologies with HPV integration analysis in a multiplex format, shows the potential to identify HPV16 integration junctions in series of clinical samples in parallel and at the same time provides E1-E2 sequences suitable for mutation/variant analysis.

aa	amino acid
ASP16	amplification selection pyrosequencing of human papillomavirus type 16
bp	base pair
CA	cancer
CF	cutoff
CGE	Cancéropôle du Grand-Est
chrom.	chromosome
CIN	cervical intraepithelial neoplasia
DKFZ	Deutsches Krebsforschungszentrum (German Cancer Research Center)
emPCR	emulsion polymerase chain reaction
GPIUA	GenomePlex universal adapter
HPV	human papillomavirus
hr-HPV	high-risk human papillomavirus
HSIL	high-grade squamous intraepithelial lesion
kb	kilo base pair
lr-HPV	low-risk human papillomavirus
LSIL	low-grade squamous intraepithelial lesion
MDA	multiple displacement amplification
nt	nucleotide
ORF	open reading frame
pBS	pBluescript vector
PCR	polymerase chain reaction
pos.	position
qPCR	quantitative polymerase chain reaction
RT-qPCR	real-time quantitative polymerase chain reaction
RA	Roche-A
RB	Roche-B
RCA	rolling circle amplification
RS-PCR	restriction-site polymerase chain reaction
TNFalpha	tumor-necrosis-factor alpha
URR	upstream regulatory region
WGA	whole genome amplification

Table of contents

1. INTRODUCTION.....	1
1.1 HUMAN PAPILLOMAVIRUS AND CERVICAL CANCER	1
1.2 HPV GENOME ORGANIZATION	2
1.3 HPV INTEGRATION AND CERVICAL CARCINOGENESIS	5
1.4 IDENTIFICATION OF HPV INTEGRATION SITES WITH THE NOVEL ASP16 STRATEGY	8
1.5 CGE-DKFZ COLLABORATION PROGRAM	10
1.6 GOALS OF THIS WORK	11
2. RESULTS.....	13
2.1 ANALYSIS OF HPV68-POSITIVE CELL LINES AND CERVICAL SCRAPES	13
2.1.1 Determination of the complete integrated HPV68b sequence in ME180.....	13
2.1.2 Determination of the complete integrated HPV68b sequence in ME180R.....	19
2.1.3 Isolation and analysis of TNFalpha-resistant cells from ME180	23
2.1.4 Cloning and sequencing of a complete HPV68b genome from a CIN2 lesion .	25
2.1.5 HPV68 variant analysis based on URR region.....	34
2.2 ANALYSIS OF HPV16 SEQUENCES GENERATED BY THE ASP16 STRATEGY	38
2.2.1 Development of computer programs for ASP16 analysis.....	38
2.2.1.1 Designing computer program strategies and platforms	38
2.2.1.2 Computer program algorithms and their tasks	41
2.2.1.3 Requirements prior to executing the programs	51
2.2.1.4 Executing the programs	53
2.2.1.5 Output files used for ASP16 analysis.....	54
2.2.2 Analysis of HPV16 sequences in ASP16-3 and ASP16-4	55
2.2.2.1 Optimization of ASP16.....	56
2.2.2.2 DNA samples in ASP16-3 and ASP16-4	58
2.2.2.3 HPV16 amplicon preparation for ASP16-3 and ASP16-4 sequencing.....	58
2.2.2.4 Amplicon sequencing by Roche/454 GS-FLX pyrosequencing	62
2.2.2.5 Statistics of ASP16-3 and ASP16-4.....	62
2.2.2.6 HPV16 integration junctions.....	67
2.2.2.7 Locations of HPV16 integration junctions.....	89

2.2.2.8 Sequence coverage of the HPV16 E1-E2 region by ASP16 sequence reads	91
2.2.2.9 Nucleotide mutation analysis in HPV16 E1-E2 area and HPV16 variant classification	94
2.2.2.10 Summary of ASP16 results for HPV16 integration and mutation analysis.....	103
3. DISCUSSION.....	104
3.1 COMPLETE SEQUENCE OF INTEGRATED HPV68B IN ME180 AND ME180R	104
3.2 FULL-LENGTH HPV68B GENOMES	105
3.3 HPV68 SUBTYPES AND VARIANTS	106
3.4 OPTIMIZED ASP16 STRATEGY AND COMPUTER ANALYSIS PROGRAMS.....	106
3.5 POSSIBLE CONSEQUENCES OF HPV16 INTEGRATION IN SAMPLES HSIL-75857 AND CIN2/3-1801	110
4. MATERIALS AND METHODS	113
4.1 CHEMICALS, COMMERCIAL MEDIA/SOLUTIONS AND ANTIBIOTICS	113
4.2 BUFFERS, STOCK SOLUTIONS AND MEDIA.....	114
4.3 COMMERCIAL KITS AND ENZYMES	117
4.4 LABORATORY EQUIPMENTS AND COMMERCIAL MATERIALS	117
4.5 COMPUTER SOFTWARE AND INTERNET RESOURCES	118
4.6 DNA MOLECULAR MARKERS	119
4.7 OLIGONUCLEOTIDE PRIMERS	120
4.7.1 Primers for cellular genes.....	120
4.7.2 Primers for HPV68 analysis.....	120
4.7.3 Primers for HPV16 and HPV16 integration junctions.....	121
4.7.4 HPV16 primers for ASP16 strategy.....	122
4.8 CERVICAL CARCINOMA CELL LINES	124
4.9 CLINICAL DNA SAMPLES.....	124
4.10 <i>IN VITRO</i> CULTIVATION OF CERVICAL CARCINOMA CELL LINES	125
4.11 ISOLATION OF TNFALPHA-RESISTANT CELLS FROM ME180	125
4.12 TNFALPHA CYTOTOXICITY ASSAY.....	126
4.13 ISOLATION OF GENOMIC DNA FROM MONOLAYER CELL CULTURES.....	126
4.14 POLYMERASE CHAIN REACTION (PCR)	127

4.15 PURIFICATION OF PCR PRODUCTS	129
4.16 CLONING OF PCR PRODUCTS	130
4.17 PLASMID DNA PREPARATION	130
4.18 DNA SEQUENCING	131
4.19 SOUTHERN HYBRIDIZATION.....	131
4.20 DNA AMPLIFICATION WITH PHI29 DNA POLYMERASE	133
4.21 PREPARATION OF PBLUESCRIPT (PBS) VECTOR FOR CLONING AT ECORI.....	134
4.22 ASP16 STRATEGY FOR HPV16 INTEGRATION ANALYSIS.....	135
4.22.1 <i>GenomePlex whole genome amplification</i>	135
4.22.2 <i>HPV16 enrichment from GenomePlex DNA libraries</i>	135
4.22.3 <i>Size selection of HPV16 multiplex PCR products</i>	137
4.22.4 <i>Roche/454 GS-FLX pyrosequencing</i>	138
4.22.5 <i>Data analysis</i>	138
REFERENCES.....	139
APPENDIX.....	148
A1. NUCLEOTIDE SEQUENCES OF INTEGRATED HPV68B IN ME180 AND ME180R	148
A2. NUCLEOTIDE SEQUENCES OF HPV68B GENOMES IN CIN2 SAMPLE	151
A3. PARTIAL URR SEQUENCES OF ELEVEN HPV68 POSITIVE SAMPLES	153
A4. NUCLEOTIDE SEQUENCES OF IDENTIFIED HPV16 INTEGRATION JUNCTIONS IN ASP16 EXPERIMENTS	154
A5. REPRESENTATIVE HPV16 E1/E2 SEQUENCES OF 25 DNA SAMPLES ANALYZED IN ASP16 EXPERIMENTS	155
A6. NUCLEOTIDE SEQUENCES OF HPV16 ORF E6 OF 25 DNA SAMPLES IN ASP16 EXPERIMENTS	167
A7. CONTENTS OF THE FILE REQUIRED FOR ASP16 DATA ANALYSIS	169
A8. SOURCE CODES OF ASP16 DATA ANALYSIS COMPUTER PROGRAMS	170

1. Introduction

1.1 Human papillomavirus and cervical cancer

Cervical cancer is the third most frequent cancer in women worldwide; with the current number of new cases annually over 500,000 and the number of death over 270,000 (de Sanjose et al, 2010; WHO/ICO Information Centre on HPV and Cervical Cancer (HPV Information Centre), 2010). Human papillomaviruses (HPVs) are detected in 99.7% of cervical cancer worldwide (Walboomers et al, 1999). It was established that infection with HPV is a necessary cause of cervical cancer (Trottier & Franco, 2006; Walboomers et al, 1999; zur Hausen, 2002), and Prof. Harald zur Hausen was awarded the 2008 Nobel Prize in Physiology and Medicine for his discovery of HPV causing cervical cancer. HPVs are classified into “types”, “subtypes” and “variants” based on the difference in nucleotide sequence of the L1 gene, with at least 10%, between 2-10% and maximally 2%, respectively (de Villiers et al, 2004). There are over 100 HPV types identified (Bernard et al, 2010; de Villiers et al, 2004), infecting cutaneous epithelia and mucosal (anogenital) epithelia. The genital HPVs, which cause benign and malignant lesions in the anogenital tract, are classified according to their potential to induce malignant transformation into high-risk and low-risk types (hr-HPV and lr-HPV). Fifteen HPV types are considered as high-risk: 16, 18, 31, 33, 35, 39, 45, 51, 52, 56, 58, 59, 68, 73 and 82 (Jacobs et al, 1995; Munoz et al, 2003). HPV16 and HPV18 are the two most prevalent genital HPVs, detected in ~50% and ~15% of cervical cancer worldwide, respectively (Castellsague, 2008; Li et al, 2010).

There are two major cervical cancer types, squamous cell carcinomas (SSCs) and adenocarcinomas. SSCs, accounting for almost 90% of cervical carcinomas (Vizcaino et al, 1998), develop from non-invasive precursor lesions, called cervical intraepithelial neoplasias (CINs) or squamous intraepithelial lesions (SILs) (Schiffman et al, 2007; Snijders et al, 2006). CINs, classified based on tissue histology, are divided into three grades, CIN1-3. SILs are classified into low and high grade, LSIL and HSIL, according to the Bethesda classification (Solomon et al, 2002). Figure 1.1 illustrates the different

2 Introduction

developmental stages of squamous epithelium from normal to cancer, in combination with CINs and SILs grades.

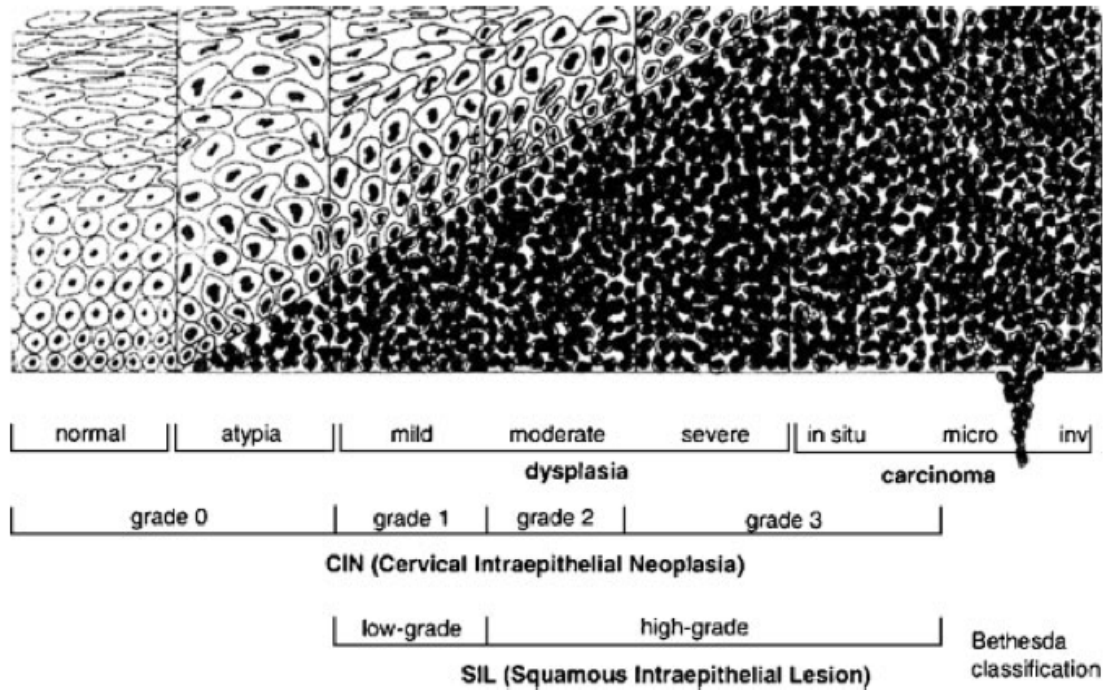


Figure 1.1: Model of histologically classified cervical squamous lesions. Taken from (Snijders et al, 2006). Morphological alterations of different development stages of cervical epithelial lesions are shown in comparisons with CIN and SIL classifications.

1.2 HPV genome organization

HPVs are small, non-enveloped DNA viruses. Their genomes are double-stranded DNA of about 8 kb encoding 8 genes, which are transcribed as polycistronic mRNAs from one DNA strand (Figure 1.2). These circular genomes are also called episomes. The HPV genome is composed of three regions: the non-coding upstream regulatory region (URR), and the protein-coding early and late regions.

The URR contains the origin of viral DNA replication, enhancer elements with binding sites for many cellular and viral transcriptional activators and repressors, and the early promoter (P97 in HPV16) at which transcription of the early genes is initiated (Park et al, 1995). The HPV16 P97 promoter (equivalent to HPV31 P99 and HPV18 P105), located upstream of ORF E6, is responsible for the expression of almost all early genes, whereas

the HPV16 P670 promoter, located within ORF E7, is responsible for the late gene expression (Zheng & Baker, 2006). The early promoter is controlled primarily by cis-elements in the URR.

The early region contains six early genes E6, E7, E1, E2, E4 and E5. The E1 and E2 genes encode proteins involved in the initiation of viral DNA replication (McBride, 2008). E1 encodes the primary viral DNA replication initiator protein, consisting of an N-terminal domain, a DNA-binding domain, an oligomerization domain and a C-terminal helicase domain (McBride, 2008; Stenlund, 2003). In the process of initiation of viral DNA replication, a dimer of the E1 protein cooperates with a dimer of the E2 protein and bind to the adjacent E1 and E2 binding sites (E1BS and E2BS) on the replication origin in the URR. The complex with E2 allows E1 to bind to the origin with higher specificity (Stenlund, 2003). When the E2 protein dimer is cut off from the complex through ATP hydrolysis by E1, additional E1 molecules are incorporated into the existing E1 dimer to form double E1 hexamers (Stenlund, 2003). The E1 hexamers melt and unwind the origin through ATP-dependent helicase activity, and when associated with the cellular DNA synthesis machinery, such as DNA polymerase alpha-primase, new viral DNA can be synthesized (McBride, 2008; Stenlund, 2003). The E2 multi-functional protein consists of three domains: the N-terminal transactivation domain (TAD), the non-conserved hinge region, and the C-terminal DNA binding domain (Ham et al, 1991). Apart from its involvement in the viral DNA replication initiation, E2 protein plays an important role in the regulation of viral promoter activity through several E2BSs in the URR (McBride, 2008). It had been demonstrated that the E2 of the bovine papillomavirus type 1 (BPV1) can activate early gene transcription (Spalholz et al, 1985; Thierry & Yaniv, 1987). When E2 binds to the E2BSs overlapping key promoter elements, such as TATA box and SP1 binding site, it represses viral transcription (McBride, 2008; Thierry, 2009). Another essential function of E2 is partitioning and maintaining the replicating viral episomes. The E2 protein binds to E2BSs in the URR of the episomes and attaches them to mitotic chromosomes, ensuring that the episomes are maintained within the nuclear envelope after mitotic and segregated evenly between daughter cells (McBride et al, 2006).

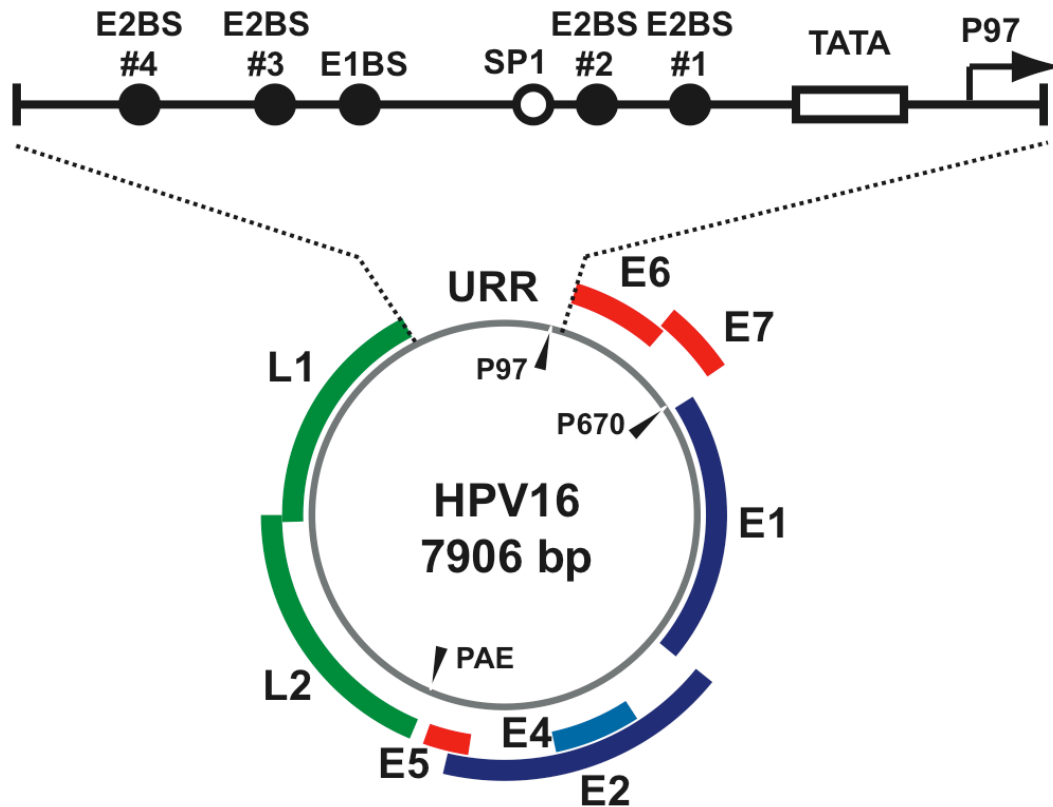


Figure 1.2: HPV16 genome organization. Adapted from Fig. 2 of (Doorbar, 2006). HPV16 genome is shown in circular form. The open reading frames (start to stop codon) of early genes (red and blue bars) and late genes (green bars) are indicated. Locations of the E1 binding site (E1BS), four E2 binding sites (E2BS #1-4), the SP1 binding site (SP1), TATA box and P97 in the non-coding upstream regulatory region (URR) are shown on top. P97 and P670: promoters. PAE: early polyadenylation site.

ORF E4 is located in the early region, embedded into the E2 region, but encodes a late gene product. Since ORF E4 does not contain a start codon, the E4 protein is translated from a transcript E1^{E4}. The E4 protein, also called E1^{E4}, therefore, carries five N-terminal amino acid residues from ORF E1. E4 protein is associated with the intermediate filament network of upper epithelial cells resulting in keratin reorganization and keratin depletion, and it may be involved in the release of mature virions from the infected cells (Doorbar et al, 1991; McIntosh et al, 2010; Wang et al, 2004). The E6, E7 and E5 genes encode viral oncoproteins (Moody & Laimins, 2010), with E6 and E7 playing the major roles in HPV-induced carcinogenesis. E6 protein binds to and facilitates the degradation of the tumor suppressor protein p53 (Werness et al, 1990). E7 protein binds to and facilitates the degradation of the tumor suppressor retinoblastoma protein (pRb) (Boyer et al, 1996; Chellappan et al, 1992). E5 protein supports cell cycle progression, and has weak transforming activity (Fehrmann et al, 2003; Maufort et al, 2010).

The late region contains two late genes L1 and L2, encoding the viral major and minor capsid proteins which encapsidate newly synthesized viral episomes (Moody & Laimins, 2010).

1.3 HPV integration and cervical carcinogenesis

Persistent infection with hr-HPV is an essential prerequisite for cervical carcinogenesis (zur Hausen, 2002). Expression of E6 and E7 genes is necessary but not sufficient to cause cervical carcinomas, and additional viral and cellular genetic alterations are required (zur Hausen, 1999). The progression from normal histology to invasive cancer is a multi-step process, as shown in Figure 1.3 (reviewed in details by Schiffman et al, 2007; Snijders et al, 2006; zur Hausen, 2000). After infection with hr-HPV, persistent expression of viral DNA combined with other factors leads to the development of precursor lesions, LSIL and HSIL (or CIN1-3). During these stages, the majority of the lesions can spontaneously regress (Chan et al, 2003; Ho et al, 1998; Nasiell et al, 1986). In some high-grade lesions, integration of hr-HPV into the host genome takes place. The integration, combined with additional cellular genetic alterations, induces genomic instability leading to progression toward invasive cancer.

In productive HPV infection (reviewed in details by Doorbar, 2005), the HPV DNA is maintained as episome of about 50-100 copies/cell (Stanley et al, 1989) in differentiating cells arising from the basal epithelial layer. During this stage, the expression of E6/E7 oncogenes is highly regulated and does not pose any risk toward cancer progression. This is because the viral oncogenes are expressed in the cells that are migrating from the inner layer to the outer epithelium, and eventually, these cells will be shed from the body. As the lesions progress, integrated HPV DNA can be detected: 5% in HSIL and 81% in cervical cancers (Cullen et al, 1991). These numbers suggested the importance of HPV integration as a major event before progression into cancer, and the integration of hr-HPV has been shown to be an important indication of progression to invasive cancer (Hopman et al, 2004).

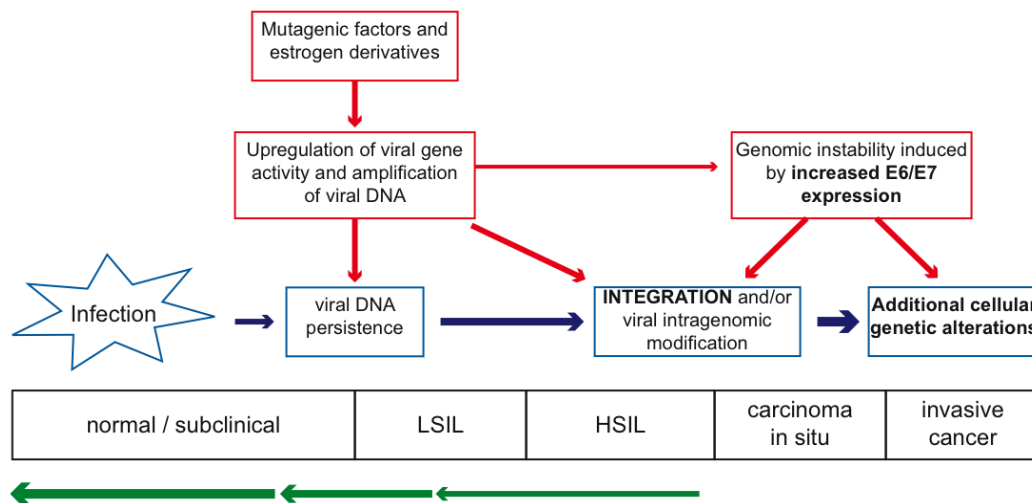


Figure 1.3: Multi-step progression to cervical cancer. Adapted from (zur Hausen, 2000). The scheme shows the progression from infection to invasive cancer. The factors contributing to progression (red and blue boxes) are shown in parallel with the histological stages (black boxes). Most of the lesions spontaneously regress back to normal as indicated by green arrows.

Upon integration, the episomal HPV DNA is disrupted and inserted into human chromosomes. The HPV breakpoints usually locate in the E1-E2 and L1-L2 regions. In integrated HPV DNA, E2 and/or E1 genes are often disrupted or deleted while URR, E6 and E7 genes always remain intact (Schwarz et al, 1985) (Figure 1.4). Since E2 represses the early promoter activity, the loss of E2 due to HPV integration allows transcriptional activation of E6/E7 oncogenes from the early promoter. HPV integration leads to increased level and stability of transcripts encoding the E6/E7 oncogenes (Jeon & Lambert, 1995), and to growth advantage of the cells (Jeon et al, 1995). Furthermore, it has been shown that disruption of E1 or E2 increases the capacity of HPV16 to immortalize cells (Romanczuk & Howley, 1992). Since the viral polyadenylation signal sequence (PAE) is missing due to integration, transcription of the E6/E7 oncogenes from the integrated HPV DNA continues until the next polyadenylation signal in the cellular DNA. This results in viral-cellular fusion transcripts. From the spliced viral-cellular fusion transcripts, the E6/E7 oncoproteins are expressed (Schneider-Gadicke & Schwarz, 1986; Smotkin & Wettstein, 1986).

HPV integration may lead to alterations of cellular genes through insertional mutagenesis. These alterations, especially the activation of cellular proto-oncogenes or inactivation of cellular tumor suppressor genes, could also contribute to cervical carcinogenesis. The integration of HPV DNA into human chromosomes occurs randomly, with preference for

genomic fragile sites and transcriptionally active regions (Kraus et al, 2008; Thorland et al, 2003; Wentzensen et al, 2004). Recurrent HPV integration sites have also been reported, in particular, integration within or close to the proto-oncogene c-myc (Couturier et al, 1991; Durst et al, 1987; Ferber et al, 2003; Peter et al, 2006; Wentzensen et al, 2002; Xu, 2010). MYC gene expression is activated when integrated hr-HPV DNA is located within chromosome 8q24 in the myc locus (Couturier et al, 1991; Peter et al, 2006; Xu, 2010), thus demonstrating activation of a cellular proto-oncogene by insertional mutagenesis. An example of the inactivation of a tumor suppressor gene by HPV integration is shown in the cervical cell line ME180 where hr-HPV68 is integrated into an intron of the potential tumor suppressor gene ZBTB7C (APM1) (Reuter et al, 1998).

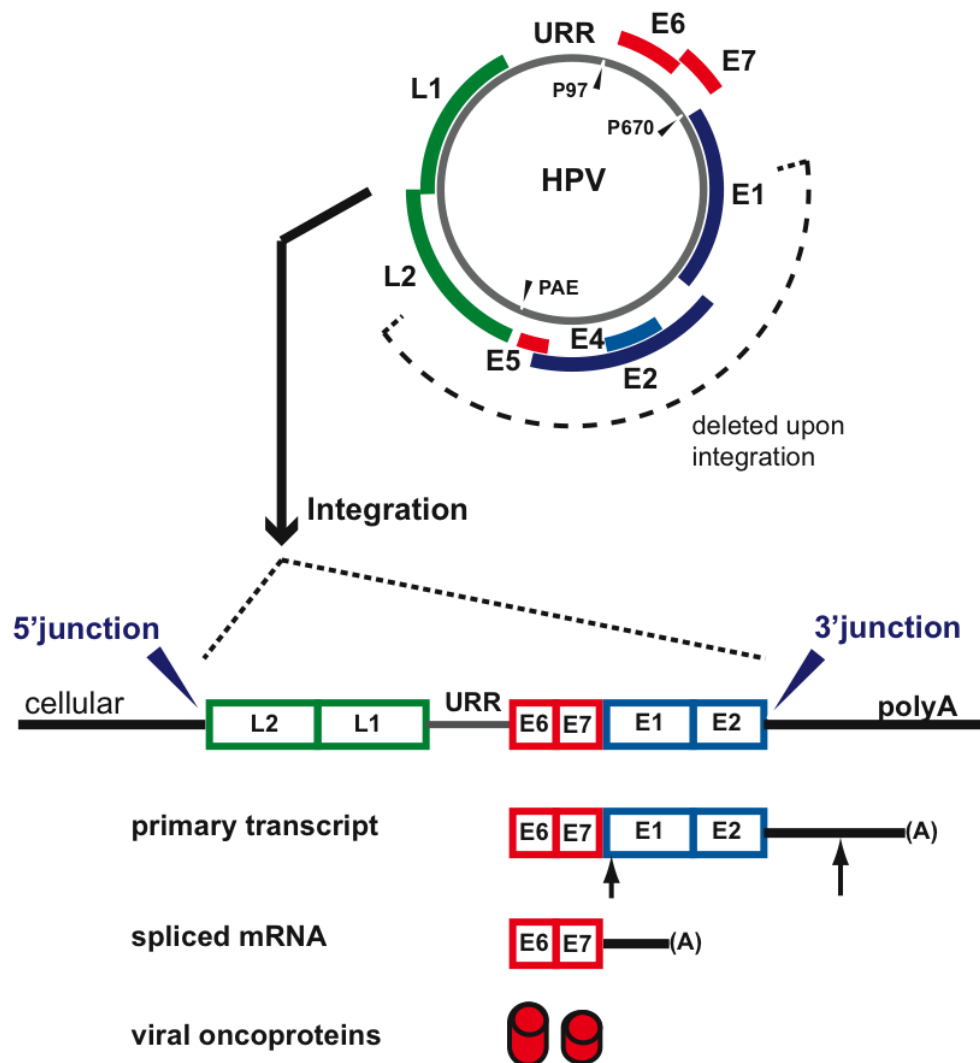


Figure 1.4: Viral oncogenes E6/E7 expression from integrated HPV DNA. An intact HPV16 episome is shown on top. The integrated DNA, flanked by cellular DNA (black lines), is shown in the lower portion. Splice donor and splice acceptor signals located in the primary transcript are indicated with small vertical black arrows. P97 and P670: early and late promoters. PAE: early polyadenylation signal.

1.4 Identification of HPV integration sites with the novel ASP16 strategy

Several methods have been developed for determining the location and structure of integrated HPV DNA. Fluorescence in situ hybridization (FISH) has been used to locate integrated HPV in the chromosomes and also to estimate the copy number of the integrated HPV (Callahan et al, 1992; Hopman et al, 2004; Mincheva et al, 1987). Based on the presence of viral-cellular transcripts produced from integrated HPV (Schneider-Gadicke & Schwarz, 1986), an RT-PCR method, named “amplification of papillomavirus oncogenes transcripts” (APOT) has been established which allows identification of integrated HPV and the chromosomal location from the cDNA sequences of the hybrid mRNA transcripts (Klaes et al, 1999). PCR-based methods for genomic DNA analysis were also developed, allowing the exact sequences of the HPV integration junctions to be determined (Kalantari et al, 2001; Luft et al, 2001). These methods require extensive laboratory works and the analysis is performed sample by sample.

In our group, another PhD student, Bo Xu, had started developing a novel PCR-based strategy for simultaneous determination of the 3' junction of integrated HPV (see Figure 1.4) of multiple DNA samples in one experiment. The strategy is called “Amplification Selection Pyrosequencing of HPV16”, abbreviated ASP16 (Xu, 2010). The principle of the ASP16 analysis includes four main steps: (1) amplification of the genomic DNA from HPV16-positive cervical lesions using GenomePlex whole genome amplification, (2) enrichment of HPV16 DNA by linear amplification and HPV16-specific multiplex PCR, (3) Roche/454 GS-FLX pyrosequencing, and (4) data analysis. The strategy is described in details in (Xu, 2010), and illustrated in Figure 1.5. Briefly, 10-30 ng total DNA of each DNA sample is amplified using GenomePlex Whole Genome Amplification (Sigma Aldrich) where the genomic DNAs are chemically fragmented, ligated with the GenomePlex universal adapter (GPUA) at both ends, and amplified by PCR. The amplified DNA product of each DNA sample, called OmniPlex library, is composed of amplified randomly fragmented total DNA with the GPUA attached at both ends. To enrich HPV16 DNA, a set of biotin-labelled HPV16-specific forward primers is used to linearly amplify HPV16 DNA-containing molecules from each OmniPlex library in two multiplex reactions, each reaction composed of eight HPV16 primers. These primers span the HPV16 E1-E2 area where the 3' junction between viral and cellular DNA is located

(see Figure 1.4). The linear amplification products are purified using streptavidin-coated magnetic beads. They are then used as templates for semi-nested HPV16 multiplex PCR using corresponding bipartite nested forward primer and a tripartite reverse primer. The forward primers are composed of the Roche-A sequence and different HPV16 sequences. The reverse primers contain the Roche-B sequence, followed by a 4-nt barcode sequence used as identification tag for each DNA sample, and the GPU A sequence (see Figure 1.5). Roche-A and Roche-B are adapter sequences required in the sequencing step. The HPV16 multiplex PCR products are sequenced using the next generation sequencing Roche/454 GS-FLX pyrosequencing system, where individual amplicon molecules are amplified by emulsion PCR (emPCR) before being sequenced (<http://www.454.com/>) (Metzker, 2010; Shendure & Ji, 2008).

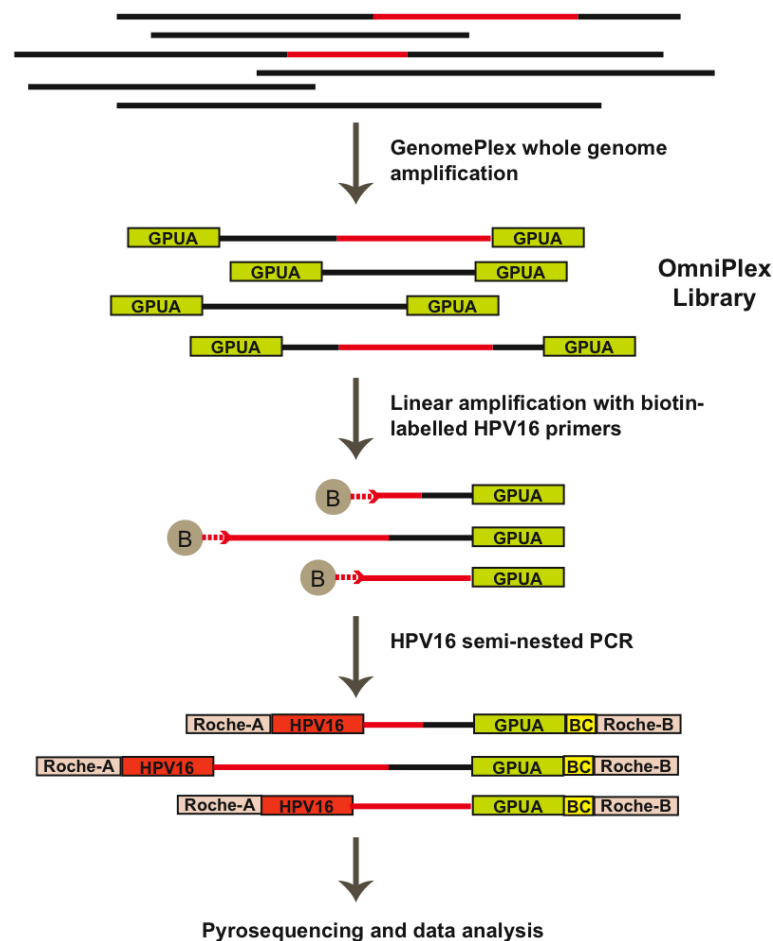


Figure 1.5: ASP16 strategy. Total DNA molecules are shown at the top. Black and red lines represent cellular DNA and HPV16 DNA, respectively. Red dashed arrows represent the biotin-labelled (B) HPV16 primers. Rectangular boxes represent the bipartite nested forward primers (Roche-A sequence fused with different HPV16 primer sequences) and the tripartite reverse primers (Roche-B sequence fused with 4-nt unique barcode sequence and GPU A sequence). GPU A: GenomePlex universal adapter. Roche-A and Roche-B: Roche adapter sequences. HPV16-boxes represent HPV16 nested primer sequences. BC: barcode.

1.5 CGE-DKFZ collaboration program

This PhD work is part of the joint program “Human papillomaviruses and cancer” between Cancéropôle du Grand-Est (CGE) and DKFZ. The main goals of the program are to unravel mechanisms of HPV-induced carcinogenesis, to identify novel progression markers, and to devise innovative therapies for virus-related cancers. We collaborate with CGE groups in Reims and Besancon working on HPV integration analysis in clinical DNA samples.

In Reims and Besancon, clinical cervical cell samples are collected from women participating in routine cervical cancer screening programs. The collected cells are tested cytologically using liquid-based ThinPrep system (Bory et al, 2002; Briolat et al, 2007; Clavel et al, 2001; Saunier et al, 2008). In the ThinPrep system, cervical cells are collected with cytobrushes, then resuspended and fixed in PreservCyt® solution. Cells are examined cytologically and classified based on Bethesda classification (Clavel et al, 2001; Solomon et al, 2002), see Figure 1.1. To determine the presence of hr-HPVs (types 16, 18, 31, 33, 35, 39, 45, 51, 52, 56, 58, 59 and 68), each cell sample is tested with the commercial Hybrid Capture 2 (HC-II) test (Digene). In the samples with hr-HPV, the HPV genotypes are determined using the commercial Linear Array (LA) HPV genotyping test (Roche) which tests for 15 hr-HPVs, 3 potentially hr-HPVs, 15 lr-HPVs and 4 HPVs with undetermined risk (Briolat et al, 2007). Additionally, cell samples are examined for DNA aneuploidy using DNA image cytometry because DNA aneuploidy has been recognized as a marker for lesions at risk for progression to cancer (Lorenzato et al, 2002; Lorenzato et al, 2001). From hr-HPV-positive cell samples, DNA and RNA are isolated for further molecular analyses.

For HPV16-positive cell samples, aliquots of the extracted DNAs are used to determine the physical status of HPV16 DNA by real-time quantitative PCR (RT-qPCR) of E2 and E6 genes (Briolat et al, 2007). The E2/E6 RT-qPCR method has been developed based on the gathered reports that the E2 hinge region is usually deleted, while the E6 gene is always intact in the integrated HPV16 (Peitsaro et al, 2002). E2 and E6 copy numbers are determined by RT-qPCR using E2- and E6-specific primer pairs and probes. The E2/E6 ratios indicate the HPV16 physical status. An E2/E6 ratio of 1 indicates pure episomal

HPV16 DNA, while values <1 indicate the presence of integrated HPV16. Figure 1.6 shows four examples of HPV16 physical status together with their expected E2/E6 ratio values and the percent integration values. Samples with pure episomal HPV16 should have 0% integration, samples with only integrated HPV16 100% integration, and samples with mixtures of episomal and integrated HPV16 DNA should have integration values of $>0\%$ to 100%. Cervical DNA samples with high percent integration values were selected for ASP16 analysis in this PhD work, because they should contain integrated HPV16 DNA.




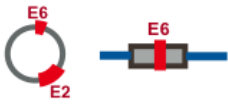
Physical status	Illustrated model	E2/E6 ratio	% Integration $= (1 - E2/E6) \times 100\%$
Episomal HPV16		$1/1 = 1$	$(1-1) \times 100\% = 0\%$
Integrated HPV16 (1 copy)		$0/1 = 0$	$(1-0) \times 100\% = 100\%$
Integrated HPV16 (2 copies)		$1/2 = 0.5$	$(1-0.5) \times 100\% = 50\%$
Mixture: 1 episome 1 integrate		$1/2 = 0.5$	$(1-0.5) \times 100\% = 50\%$

Figure 1.6: HPV physical status and E2/E6 ratio determined by RT-qPCR. The copy numbers of E6 and E2 genes (red boxes) are determined by RT-qPCR (Peitsaro et al, 2002). The E2/E6 ratios are indicated. Gray boxes represent individual HPV16 copies. Blue lines represent cellular DNA.

1.6 Goals of this work

This PhD work covers the integration analysis of two different hr-HPV types, HPV16 and HPV68b. While HPV16 is the most frequently detected type in cervical cancer, HPV68b is a rare hr-HPV. Previously in our group, HPV68 subtype b (HPV68b) was identified as integrated DNA in the cervical cell line ME180 (Reuter et al, 1991). The reported sequence of this integrated HPV68b was determined based on a cloned restriction fragment. Further analyses of this cell line, however, showed evidence of a probable

cloning-induced deletion in this restriction fragment (Reuter et al, 1991). Therefore, in the beginning of this PhD work, the complete and correct sequence of the integrated HPV68b DNA in the cell line ME180 was determined. Additionally, the mutant cell line ME180R, which is resistant to growth-inhibition by tumor-necrosis-factor-alpha (TNFalpha), was found to contain largely altered integrated HPV68b DNA compared to the parental cells ME180 (Elisabeth Schwarz, unpublished data). The question arose whether these alterations contribute to the resistance to TNFalpha. Thus, in this work, the complete integrated HPV68b DNA sequence in the cell line ME180R was determined, and it was investigated whether alterations in the integrated HPV68b are associated with resistance to TNFalpha by selecting new TNFalpha-resistant ME180 variants. Since no full-length HPV68b genome sequence was available, a complete HPV68b genome was isolated from a CIN2 DNA sample, and was sequenced.

In the HPV16 analysis using ASP16 strategy, the Roche/454 GS-FLX pyrosequencing system has the capacity to deliver about 210,000 sequence reads per sequencing round for the selected format. Such massive amount of data requires computer analysis software for data management and interpretation. With my knowledge in computer programming, I collaborated with Bo Xu in the ASP16 project and developed computer programs for automatically processing and analyzing ASP16 sequence data. Two ASP16 experiments had been performed previously (Xu, 2010). Because the strategy required further optimization, additional ASP16 experiments were performed as part of this PhD work.

2. Results

2.1 Analysis of HPV68-positive cell lines and cervical scrapes

In this study, the complete sequences of the integrated HPV68b in cell lines ME180 and ME180R were determined, and a complete genome of HPV68b was cloned from a CIN2 lesion and sequenced. By determination of sequence polymorphisms in the non-coding upstream regulatory region (URR) of HPV68 DNA in different cervical scrapes, new HPV68 variants were identified.

2.1.1 Determination of the complete integrated HPV68b sequence in ME180

The cell line ME180 was established from an omental metastasis of a rapidly spreading cervical carcinoma (Sykes et al, 1970). ME180 cells contain HPV68b DNA, integrated into chromosome 18q21 (Reuter et al, 1991). The integrated HPV68b in ME180 was first cloned and sequenced by Reuter et al (Reuter et al, 1991). The clone contains a *SacI* restriction fragment of 13.1 kb as insert AA13.1. Sequence analysis showed that the 13.1 kb fragment contains two incomplete copies of HPV68b DNA, named 5'copy and 3'copy, flanked by cellular sequences (Figure 2.1). Southern hybridization of *SacI*-digested genomic DNA showed a 20-kb restriction fragment instead of 13 kb (Reuter, 1995). This size discrepancy indicated a probable cloning-induced deletion of about 7 kb. According to the sequence organization, it was suspected that this deletion occurred in the area of the 5' copy of the integrated HPV68b (Figure 2.1). In the previous analyses, it was shown by Southern and PCR experiments of ME180 DNA that the full 5'copy probably contains ORFs E5, L2, L1, the non-coding URR, ORFs E6, E7 and E1 before the breakpoint in ORF E2 is reached (Reuter, 1995; Marco Springer and Elisabeth Schwarz, unpublished data). However, the exact structure and nucleotide sequences had not been determined. This was performed as part of this PhD thesis.

To obtain a complete sequence of the integrated HPV68b in ME180, both 5'copy and 3'copy were amplified by PCR, cloned into plasmid vectors and sequenced. Roche's Expand Long Template PCR System was used for amplification because of its high fidelity enzyme and long-range amplification capability. The amplification strategy for

both copies is illustrated in Figure 2.2. The 7.2-kb amplified fragment of the 5'copy was cloned into vector pCR4-TOPO (Invitrogen, Karlsruhe), the 6.2-kb fragment of the 3'copy into vector pSC-A (Stratagene, USA). The two clones were sequenced in both strands with Big-Dye terminator chemistry on an ABI model sequencer (Andreas Hunziker at the DKFZ Core Facilities). The sequences obtained from both clones (7123 bp and 6196 bp) were assembled at the overlapping area.

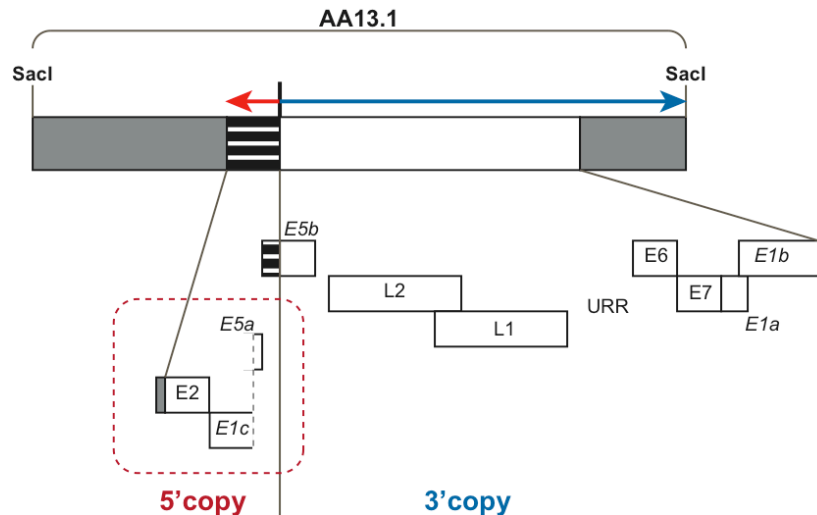


Figure 2.1: Previously cloned and sequenced *SacI* fragment AA13.1 containing integrated HPV68b in ME180. The 13.1-kb fragment contains two integrated incomplete copies of HPV68b in head-to-head orientation. The 5'copy is indicated by a horizontally striped box, the 3'copy by a white box, and the flanking cellular sequences by gray boxes. Red and blue arrows indicate the 5' to 3' polarity of the coding DNA strand. The open reading frames (ORFs) are shown in the enlarged segment below. The 5'copy was suspected to contain a cloning-induced deletion (circled in dashed red line). This figure is adapted from Fig. 2 of Reuter et al (Reuter et al, 1991).

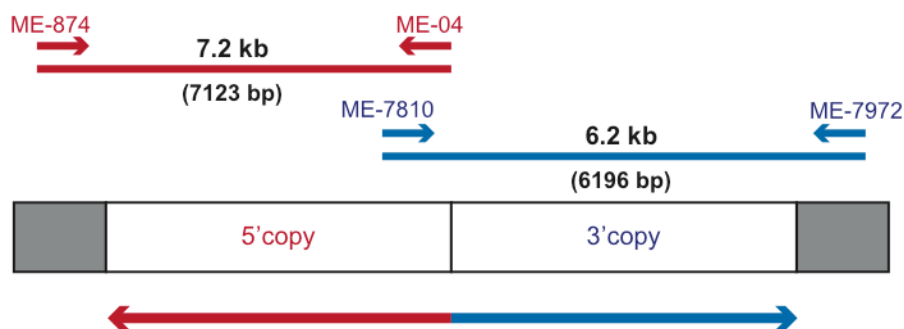


Figure 2.2: Amplification of the integrated HPV68b in ME180. Two fragments (red and blue lines), covering the integrated 5'copy (long red arrow) and 3'copy (long blue arrow) of HPV68b DNA were amplified from the genomic DNA of cell line ME180. The two primer pairs are indicated (short red and blue arrows) as well as the exact amplicon sizes determined by sequencing. Primers ME-874 and ME-7972 bind to the cellular DNA of chromosome 18 upstream and downstream of the HPV68b integration site. Primers ME-04 and ME-7810 bind to HPV68b DNA. Primer sequences are indicated in Materials and Methods.

The assembled nucleotide sequence, named HPV68b(int)-ME180, contains 13286 nucleotides. The sequence is shown in Appendix A1. The genetic organization of HPV68b(int)-ME180 is shown in Figure 2.3 and Table 2.1, and the following characteristics are emphasized:

1. The two HPV68b copies are integrated in head-to-head orientation into chromosome 18 between positions 45578341 (5' junction) and 45578364 (3' junction) based on sequence accession NC_000018.9, with 22 bp of cellular DNA deleted.
2. The 5' copy contains 6993 bp of HPV68b DNA. The 5' copy in AA13.1 has a size of only 870 bp (Reuter et al, 1991). This confirms the previous assumption of a cloning-induced deletion. Viewed in the 5'-to-3' direction of the coding strand, the 5' copy contains 41-bp non-coding sequence, intact ORFs E5, L2 and L1, the 809-bp URR, intact ORFs E6 and E7, followed by incomplete ORF E1 and truncated ORF E2. The identification of the 41-bp sequence as HPV68b DNA was possible because the sequence of a complete HPV68b genome has also been determined in this PhD work (see section 2.1.4). This sequence is non-coding, located between the E2 stop codon and the E5 start codon. The ORF E5 was considered to be complete because it contains the same 222 nucleotides from the start codon to the stop codon as in the published ORF E5 of HPV68a (accession DQ080079). The URR contains a 2-bp gap compared to the 3' copy (Figure 2.19). The ORF E1 harbors a 64-bp deletion causing a frameshift and a disruption into ORFs E1a and E1b (Figure 2.4). The E2 gene carries the integration breakpoint located 359 bp downstream of the E2 ATG start codon.
3. The 3' copy contains 6071 bp of HPV68b DNA. Viewed in the 5'-to-3' direction of coding strand, the 3' copy contains a truncated ORF E5, intact ORFs L2 and L1, intact URR, followed by intact ORFs E6 and E7, and a truncated ORF E1. Compared to the complete E5 ORF in the 5' copy, the E5 ORF of the 3' copy is truncated, missing the first 46 bp. The E1 gene is truncated at 1289 bp downstream of the ATG start codon. It is joined with a 14-bp sequence derived from ORF E5 before the flanking cellular sequences start.
4. The two viral-cellular junctions in the 5' and 3' copies of HPV68b(int)-ME180 are identical to those of AA13.1 (Figure 2.3).
5. No intact ORF E1 is present in HPV68b(int)-ME180. The E1 of the 3' copy is truncated due to integration and the E1 of the 5' copy harbors a 64-bp deletion and frameshift (Figure 2.4). Surprisingly, the 64-bp deletion originally identified in AA13.1 in

the E1 part of the 3'copy is present in the complete sequence HPV68b(int)-ME180 in the E1 part of the 5'copy (Figure 2.3). This indicates complex cloning-induced DNA rearrangements in clone AA13.1. In the 5'copy, ORF E1a encodes 116 amino acids with 108 N-terminal translated from the E1 frame and 8 C-terminal amino acids from a shifted frame due to the 64-bp deletion. ORF E1b encodes 555 amino acids, with 43 N-terminal residues translated from a shifted frame of E1 due to the 64-nt deletion and 512 C-terminal residues from the normal frame. In the 3'copy, the truncated ORF E1 encodes 445 amino acids, with the 430 N-terminal residues translated from E1 sequence, 4 residues from the partial E5 sequence, and 11 C-terminal residues from cellular DNA. The sequence of the complete ORF E1 of HPV68b-ME180 was assembled through combination of the E1 sequences in 5' and 3' copies. The assembled intact ORF E1 contains 1923 bp encoding 640 amino acids, the same numbers as ORF E1 of the published HPV68a (DQ080079).

6. The HPV68b(int)-ME180 contains no intact ORF E2. The ORF E2 in the 5'copy is truncated due to integration and E2 is completely deleted in the 3'copy. The ORF E2 of the 5'copy encodes 134 amino acids, with 120 N-terminal residues translated from the E2 gene and 14 C-terminal residues from cellular DNA. By comparison with the complete ORF E2 of HPV68a (DQ080079), 754 bp are missing.

7. The URR and ORFs E6 and E7 are present in intact form in both 5' and 3' copies.

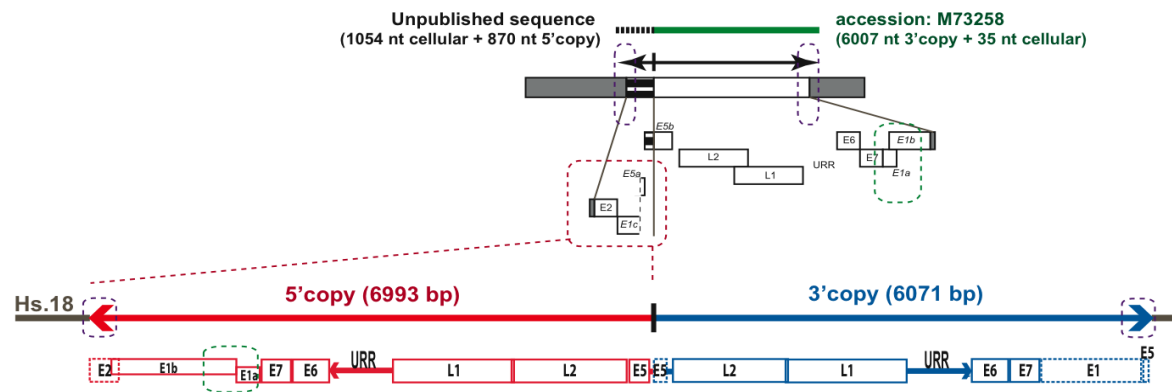


Figure 2.3: Complete structure of HPV68b(int)-ME180. The previously sequenced clone AA13.1 (Reuter et al, 1991) is shown at the top for comparison. The top green line represents the published sequence (accession M73258) containing the 3'copy of AA13.1, whereas the dashed black line represents the sequence of the 5'copy (Elisabeth Schwarz, unpublished). The locations of ORFs in the integrated HPV68b are shown at the bottom. Solid-lined red and blue boxes indicate the complete coding region (from ATG to stop codon) of each ORF, the dashed red and blue boxes the truncated ORFs. The red and blue arrows indicate the direction of transcription for the 5'copy and 3'copy, respectively. The viral-cellular junctions in the 5' and 3' copies in clone AA13.1 and in the new complete structure are identical (dashed purple circles). The E1a-E1b area of AA13.1 is present in the 5'copy of the complete structure (dashed green circle). AA13.1 contains 6007 bp HPV68b DNA in the 3'copy due to the presence of the 64-bp deletion. The positions of all ORFs are given in Table 2.1.

Table 2.1: Sequence organization of HPV68b(int)-ME180

Position	Length (bp)	Description ^(a)	Remark	HPV68b-CIN2 ^(b)
1-186	186	flanking cellular DNA of chrom.18, pos. 45578156-45578341 (NC_000018.9)	cellular DNA	
182-186		cellular-viral junction (5 bp overlap)		
136-540	405	truncated ORF E2 ^(c)	5'copy (6993 bp)	3031-2673
467-2002/2134	1536/1668	ORF E1b ^(c)		2746-1211
2001/2002		borders of 64-bp deletion		1212/1147
1975-2325	351	ORF E1a ^(c)		1238-824, (del 1211-1148)
2332-2664	333	complete ORF E7		817-485
2672-3148	477	complete ORF E6		477-1
3149-3957	809	URR		7836-7026
3958-5475	1518	complete ORF L1		7025-5508
5456-6865	1410	complete ORF L2		5527-4118
6912-7133	222	complete ORF E5		4071-3850
7134-7174	41			3849-3809
7174/7175		junction between 5'copy and 3'copy	3'copy (6071 bp)	
7175-7354	180	truncated ORF E5		3892-4071
7401-8810	1410	complete ORF L2		4118-5527
8791-10308	1518	complete ORF L1		5508-7025
10309-11119	811	URR		7026-7836
11120-11596	477	complete ORF E6		1-477
11604-11936	333	complete ORF E7		485-817
11943-13280	1338	truncated ORF E1		824-2112
13231		3'end of truncated E1		2112
13232-13245	14	part of ORF E5		3927-3914
13245		viral-cellular junction (1 bp overlap)		
13245-13286	42	flanking cellular DNA of chrom.18, pos. 45578364-45578405 (NC_000018.9)	cellular DNA	

- (a) Because the 5'copy is in opposite orientation, position numbering of the ORFs is from stop codon to start codon. In the 3'copy, positions are opposite.
 (b) Complete HPV68b genome isolated from a CIN2 lesion, described in section 2.1.4.
 (c) See text and Figure 2.4 for explanation.

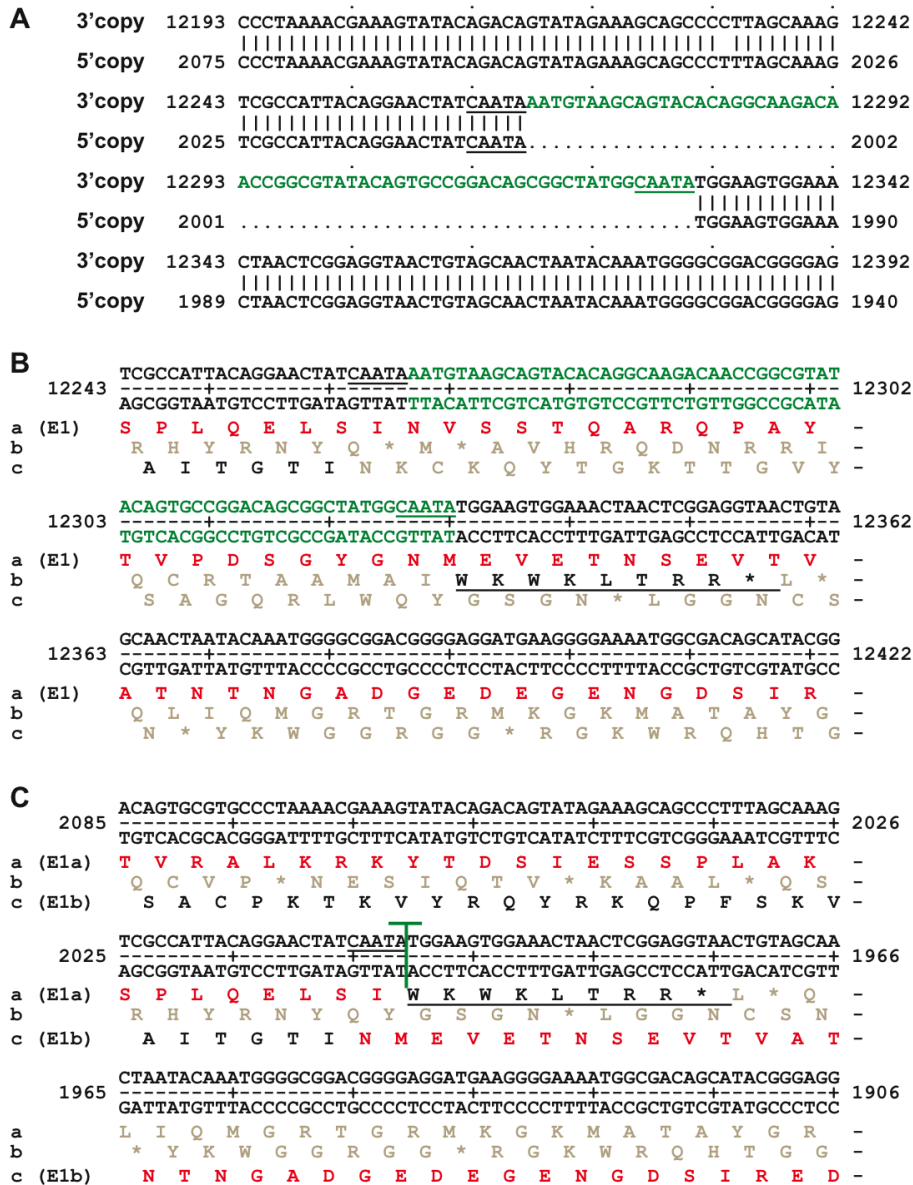


Figure 2.4: Internal 64-bp deletion and frameshift in ORF E1 of HPV68(int)-ME180. Panel A: Partial sequence comparison of E1 gene in the 5'copy and 3'copy. Green characters indicate the 64 bp deleted in the 5'copy. The repeated CAATA sequence at the border of the 64 bp deletion is underlined. **Panel B:** Partial nucleotide sequence of ORF E1 in the 3'copy, covering pos. 12243-12422. The 64-bp sequence deleted in the 5'copy is shown in green. The E1 protein is translated from frame a. The underlined amino acids in frame b indicate the frameshifted residues in ORF E1a of the 5'copy. **Panel C:** Partial nucleotide sequence of ORFs E1a and E1b in the 5'copy covering pos. 2085-1906, shown with translated amino acids. The green T-marker indicates the position of the 64-bp deletion. ORF E1a encodes 116 amino acid translated from frame-a with 8 C-terminal residues (underlined) originated from another frame in intact ORF E1 (see Panel C). ORF E1b is translated from frame-c. Red-color indicates amino acids originating from the E1 reading frame.

2.1.2 Determination of the complete integrated HPV68b sequence in ME180R

The cell line ME180R was derived from ME180 cells by selection of subclones resistant to growth inhibition by tumor-necrosis-factor-alpha (TNFalpha) (Pfreundschuh et al, 1989). Comparing the genomic organization of integrated HPV68b in ME180 and ME180R by Southern hybridization, large deletions of the integrated HPV68b DNA in ME180R became apparent (Figure 2.5; Schwarz, E, Reuter, S and Springer, M, unpublished data). Additional mapping and hybridization experiments had been performed in the past, but the exact structure of the integrated HPV68b in ME180R had not been determined. This task was performed as part of this PhD work.

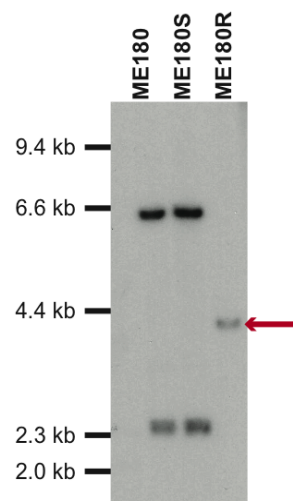


Figure 2.5: Genomic Southern hybridization of ME180, ME180R and ME180S. BamHI-digested genomic DNA from cell lines ME180, ME180R (TNFalpha resistant) and ME180S (TNFalpha sensitive) were hybridized with an HPV68b DNA probe (L1-L2-URR-E6-E7). The hybridization pattern of ME180R is different from the other two cell lines, showing a single fragment smaller than 4.2 kb (red arrow).

Based on the differences in the Southern hybridization patterns (Figure 2.5) together with the known locations of BamHI sites in the complete integrated HPV68b(int)-ME180, it was deduced that about 6200 bp of HPV68b DNA were deleted in ME180R. To determine the exact structure of the integrated HPV68b sequence in ME180R, five fragments (R_PCR_1 to R_PCR_5) were amplified using Roche's Expand Long Template PCR System. The primers were selected based on the previously determined HPV68b(int)-ME180 sequence. Table 2.2 shows the primer positions and the product sizes. The PCR products were cloned into vector pCR4-TOPO (Invitrogen, Karlsruhe), and completely sequenced. The sequences of the five fragments were assembled based on the overlapping areas (Figure 2.6).

The assembled nucleotide sequence, named HPV68b(int)-ME180R, consists of 6544 nucleotides and is shown in Appendix A1. The genetic organization of the integrated HPV68b in ME180R in comparison to ME180 is shown in Figure 2.6 and Table 2.3. Important features are explained in the following:

1. Two large deletions of HPV68b are present, one in the 5'copy (2316-bp deletion) and another in the 3'copy (4679-bp deletion). The sequences at the deletion boundaries are shown in Figure 2.7.
2. In the 5'copy, ORF L1 is completely deleted and ORF L2 partially deleted. In the 3'copy, the complete ORFs L2, L1, URR and E6, and partial E7 ORF (covering 273 bp from the ATG start codon) were deleted. The intact URR, together with the complete ORFs E6 and E7 are present in the 5'copy. Therefore, the oncogenes E6/E7 can only be expressed from the 5'copy.
3. Sequence comparison demonstrated that the region containing the head-to-head viral junction is reversed (Figure 2.6). The complete ORF E5 is present. However, in comparison to HPV68b(int)-ME180, it contains two point mutations, one of them changing a TGG into a TAG premature stop codon (Figure 2.8). Therefore, the E5 protein is shortened by 12 amino acids.
4. The URR in the 5'copy is shortened by four bp at the 5'end (Figure 2.7), and contains a 2-bp insertion (Figure 2.19), when compared with the URR of the HPV68b(int)-ME180 5'copy.

Table 2.2: Primers and products for PCR amplification of the integrated HPV68b in ME180R.

PCR name*	Primer**	5'end position in HPV68b(int)-ME180	ME180 PCR calculated product size	ME180R PCR product size***
R_PCR_1	ME-01	-247 (F)	2456 bp	~2.5 kb
	ME-6805	2209 (R) and 12059 (F)		
R_PCR_2	ME-7094	1924 (F) and 12408 (R)	5318 bp	~2.2 kb
	ME-7810	7091 (F)		
R_PCR_3	ME-1349	3775 (F) and 10491 (R)	9512 bp	~2.5 kb
	ME-7972	13286 (R)		
R_PCR_4	ME-7718	6999 (R) and 7267 (F)	5076 bp (5'copy) 5142 bp (3'copy)	~1.3 kb
	ME-7094	1924 (F) and 12408 (R)		
R_PCR_5	ME-6651	2367 (R) and 11901 (F)	1386 bp	~1.4 kb
	ME-7972	13286 (R)		

* Genomic DNA of ME180R was used as template.

** Primer sequences are given in Materials and Methods.

*** Estimated sizes from agarose gel.

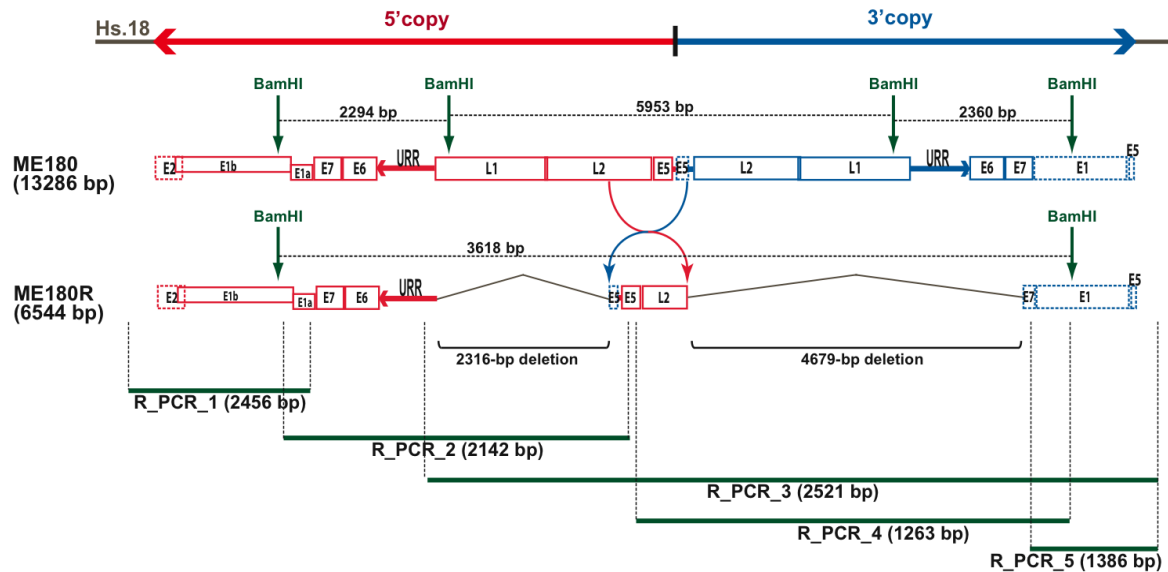


Figure 2.6: Structure of HPV68b(int)-ME180R in comparison with HPV68b(int)-ME180. The two large deletions in the 5'copy and 3'copy are shown. As indicated by the arrows, the viral-viral junction area (middle) is reversed in ME180R. BamHI restriction sites are indicated by green vertical arrows, and the sizes of the BamHI restriction fragments deduced from the nucleotide sequences are given. At the bottom part, regions covered by the five PCR products of ME180R are indicated together with the PCR sizes.

Table 2.3: Sequence organization of HPV68b(int)-ME180R in comparison with HPV68b(int)-ME180.

HPV68b(int)-ME180		HPV68b(int)-ME180R	Retained and deleted regions of ME180R*
Positions	Description ^(a)	Positions	Positions
1-186	flanking cellular DNA of chrom.18 pos. 45578156-341 (NC_000018.9)	248-433	retained: 3955 bp ME180R 248-4202 ^(b) ME180 1-3953
182-186	cellular-viral junction (5 bp overlap)	429-433	
136-540	truncated ORF E2	383-787	
467-2002/2134	ORF E1b	714-2249/2381	
2001/2002	borders of 64-bp deletion	2248/2249	
1975-2325	ORF E1a	2222-2572	
2332-2664	complete ORF E7	2579-2911	
2672-3148	complete ORF E6	2919-3395	
3149-3957	URR	3396-4202 ^(b)	
3958-5475	complete ORF L1	deleted: ME180 3954-6269 (2316 bp)	
5456-6865	complete ORF L2	5133-4538 (inversed and truncated)	retained: 928 bp ME180R 4206-5133 ME180 7197-6270
6912-7133	complete ORF E5	4491-4270 (inversed)	
7174/7175	junction between 5'copy and 3'copy	4229/4228 (inversed)	
7175-7354	truncated ORF E5	4228-4206 (inversed and truncated)	
7401-11596	complete ORFs L2, L1, URR, E6	deleted: ME180 7198-11876 (4679 bp)	
11604-11936	complete ORF E7	5135-5194 (truncated)	retained: 1410 bp ME180R 5135-6544 ME180 11877-13286
11943-13280	truncated ORF E1	5201-6538	
13231	3'end of truncated E1	6489	
13232-13245	part of ORF E5	6490-6503	
13245	viral-cellular junction (1 bp overlap)	6503	
13245-13286	flanking cellular DNA of chrom.18 pos. 45578364-405 (NC_000018.9)	6503-6544	

* Compared to HPV68b(int)-ME180.

(a) see footnote of Table 2.1.

(b) contains a 4-bp deletion at the 5' end of the URR and a 2-bp insertion (see text and Figure 2.7).

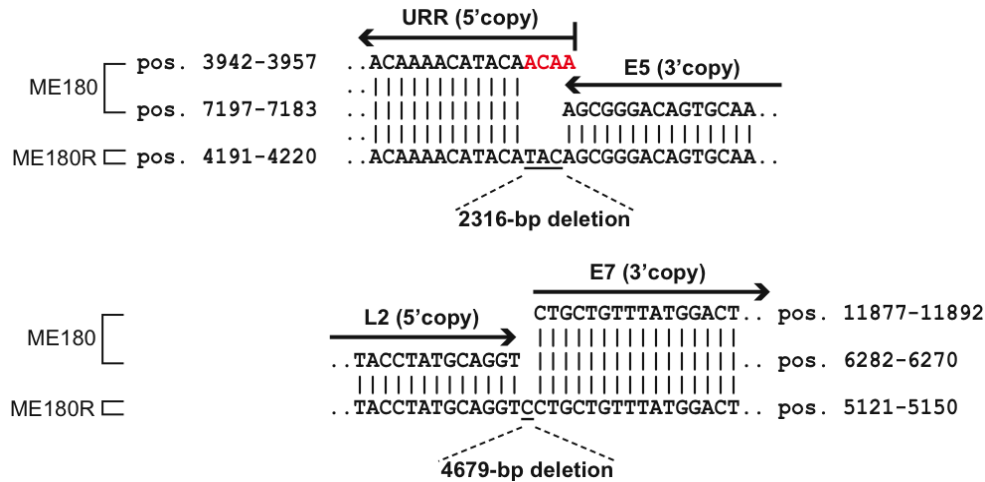


Figure 2.7: Nucleotide sequences at the deletion boundaries in ME180R. The top and bottom parts show the nucleotide sequence comparisons between the integrated HPV68b in ME180 and ME180R at the 5'copy and 3'copy deletion areas, respectively. The arrows represent the direction of transcription. The first four basepairs of the URR are deleted in ME180R (red). Three additional nucleotides (underlined TAC) are present in the 5'copy deletion area of ME180R, and a single additional nucleotide (underlined C) in the 3'copy deletion area.

A

```

ME180 7133 ATGATTGTACTGGTATTTTTGGTGTGGTTTTGTGTGTGCATGTATATATGTTGCACTGTC 7074
|||||
ME180R 4270 ATGATTGTACTGGTATTTTTGGTGTGGTTTTGTGTGTGCATGTATATATGTTGCACTGTC 4329

ME180 7073 CCGCTTCTGCAGTCCATGCATGTGTGTGTATGTGTGGATACTTGTGTTTGTGTTTATA 7014
|||||
ME180R 4330 CCGCTTCTGCAGTCCATGCATGTGTGTGTATGTGTGGATACTTGTGTTTGTGTTTATA 4389

ME180 7013 TTAGTACGTACCACACCATTTGGAGGTCTTTGCTGTATATATACTTTTTTTTTTACTGCCT 6954
|||||
ME180R 4390 TTAGTACGTACCACACCATTTGGAGGTCTTTGCTGTATATATACTTTTTTTTTTACTGCCT 4449

ME180 6953 GTGTGGGTATTACACAGTTTTGCTCGTTATAGTATGCCTTAA 6912
|||||
ME180R 4450 ATGTAGGTATTACACAGTTTTGCTCGTTATAGTATGCCTTAA 4491

```

B

```

4390 TTAGTACGTACCACACCATTTGGAGGTCTTTGCTGTATATATACTTTTTTTTTTACTGCCT 4449
+-----+-----+-----+-----+-----+-----+-----+-----+
AATCATGCATGGTGTGGTAACCTCCAGAAACGACATATATATGAAAAAATGACGGA
L V R T T P L E V F A V Y I L F F L L P -

4450 ATGTAGGTATTACACAGTTTTGCTCGTTATAGTATGCCTTAA 4491
+-----+-----+-----+-----+-----+-----+-----+-----+
TACATCCATAATGTGTCAAAACGAGCAATATCATACGGAATT
M * V L H S F A R Y S M P * -

```

C

```

7013 TTAGTACGTACCACACCATTTGGAGGTCTTTGCTGTATATATACTTTTTTTTTTACTGCCT 6954
+-----+-----+-----+-----+-----+-----+-----+-----+
AATCATGCATGGTGTGGTAACCTCCAGAAACGACATATATATGAAAAAATGACGGA
L V R T T P L E V F A V Y I L F F L L P -

6953 GTGTGGGTATTACACAGTTTTGCTCGTTATAGTATGCCTTAA 6912
+-----+-----+-----+-----+-----+-----+-----+-----+
CACACCCATAATGTGTCAAAACGAGCAATATCATACGGAATT
V W V L H S F A R Y S M P * -

```

Figure 2.8: ORF E5 with premature stop codon in ME180R. Panel A: ORF E5 of ME180R is shown in comparison with ME180 5'copy. Point mutations in ME180R are colored in blue and red. Panel B: Partial translation of ORF E5 of ME180R containing missense (blue) and nonsense (red) mutations. Panel C: Partial translation of intact ORF E5 of HPV68b(int)-ME180 5'copy.

2.1.3 Isolation and analysis of TNFalpha-resistant cells from ME180

The complete sequence of the integrated HPV68b in the TNFalpha-resistant ME180R cells revealed that the HPV68b DNA is highly altered, especially by two large deletions, in comparison to the integrated HPV68b in the parental ME180 cells. Therefore, it was the question whether these two large deletions might contribute to the TNFalpha-resistant characteristics of ME180R cells. To investigate this issue, TNFalpha-resistant cells were independently isolated from the parental ME180 cells. The selection procedure is described in Figure 2.9. Two TNFalpha-resistant variants, ME180-2A and ME180-3A, were obtained.

The two newly derived variants were tested in a cytotoxicity assay with different concentrations of TNFalpha to determine their resistance to TNFalpha, compared with ME180 and ME180R cells. The surviving cells (cell viability) of the four cell populations after TNFalpha treatment were stained with crystal violet and their absorbance measured at wavelength 560 nm (Table 2.4). The calculated percentage values of cell viabilities are shown in a graphical form in Figure 2.10. Compared to the parental ME180 cells, the ME180R and the two isolated variants ME180-2A and ME180-3A are highly resistant to TNFalpha. Both variants showed comparable resistance to TNFalpha as ME180R.

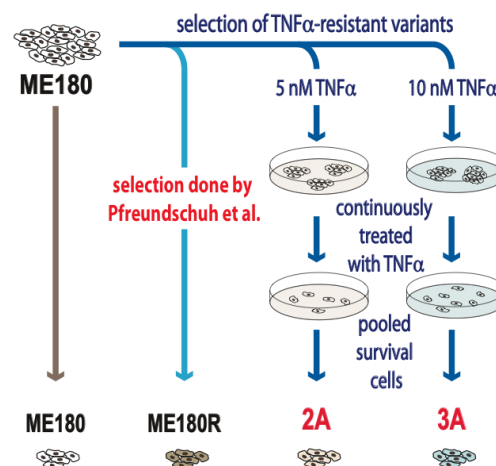


Figure 2.9: Isolation of TNFalpha-resistant variants from parental ME180 cells. ME180 cells were continuously treated with two different concentrations of TNFalpha for 18 weeks. Variants ME180-2A and ME180-3A were obtained from the pooled cells after the treatment with 5 nM and 10 nM TNFalpha, respectively. ME180R was previously selected (Pfreundschuh et al, 1989).

If the two large deletions in the integrated HPV68b in ME180R contribute to the TNFalpha-resistance phenotype, then the two newly isolated TNFalpha-resistant variants should contain some deletions as well. A genomic Southern hybridization was performed to compare the structures of integrated HPV68b in ME180, ME180R, M180-2A and ME180-3A. The results are shown in Figure 2.11. The HPV68b hybridization patterns of variants ME180-2A and ME180-3A were identical to that of ME180, demonstrating that deletions are not present. From these data, it can be concluded that the two large deletions in ME80R are not causally linked to the TNFalpha-resistant growth phenotype.

Table 2.4: Cell viability measurements for TNFalpha resistance of ME180, ME180R, ME180-2A and ME1803A.

A: Absorbance intensities of viable cells after treatment with TNFalpha. ^(a)

	ME180			ME180R			2A			3A			
0 nM	0,698	0,860	0,896	0,522	0,541	0,441	0,343	0,357	0,344	0,451	0,462	0,409	5000 cells/well
2.5 nM	0,107	0,100	0,120	0,570	0,608	0,513	0,389	0,410	0,380	0,717	0,630	0,609	
5 nM	0,107	0,111	0,109	0,574	0,622	0,506	0,422	0,409	0,430	0,742	0,649	0,685	
10 nM	0,100	0,106	0,104	0,638	0,623	0,517	0,389	0,403	0,394	0,709	0,639	0,609	
0 nM	1,018	1,021	0,990	0,583	0,581	0,477	0,390	0,394	0,372	0,703	0,578	0,485	10000 cells/well
2.5 nM	0,127	0,109	0,113	0,569	0,518	0,395	0,327	0,288	0,349	0,512	0,774	0,000	
5 nM	0,131	0,140	0,139	0,591	0,606	0,484	0,381	0,360	0,360	0,556	0,571	0,516	
10 nM	0,111	0,105	0,121	0,610	0,649	0,551	0,400	0,418	0,386	0,640	0,591	0,425	

B: Average absorbance intensities. ^(b)

[average]	ME180	ME180R	2A	3A
0 nM	0,914	0,524	0,367	0,515
2.5 nM	0,113	0,556	0,357	0,540
5 nM	0,123	0,564	0,394	0,620
10 nM	0,108	0,598	0,398	0,602

C: Percentage of absorbance intensities. ^(c)

[%]	ME180	ME180R	2A	3A
0 nM	100	100	100	100
2.5 nM	12	106	97	105
5 nM	13	108	107	120
10 nM	12	114	109	117

(a) The cells, seeded at densities of 5000 or 10000 cells/well, were treated in triplicates for each TNFalpha concentration. Gray highlighted value is excluded from the average and percentage calculation.

(b) Each value is an average of the same cell population, treated with the same TNFalpha concentration.

(c) Percentage based on cells in absense of TNFalpha (0 nM) for each cell population.

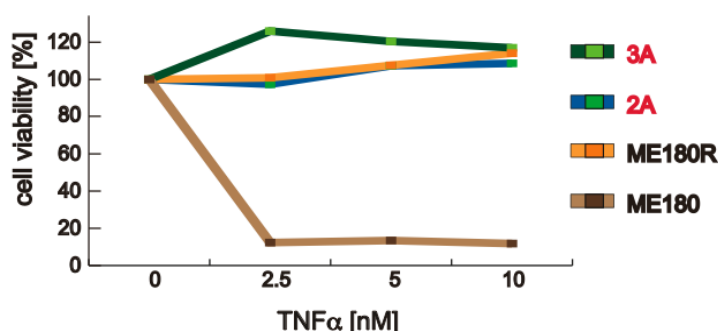


Figure 2.10: Cell viability for TNFalpha resistance of ME180, ME180R, ME180-2A and ME180-3A. The cell viabilities, determined as percentages of the number of cells after TNFalpha treatment (2.5, 5 and 10 nM) over the control (0 nM) of each cell population, are shown. The values of this plot are indicated in Table 2.4, Panel C.

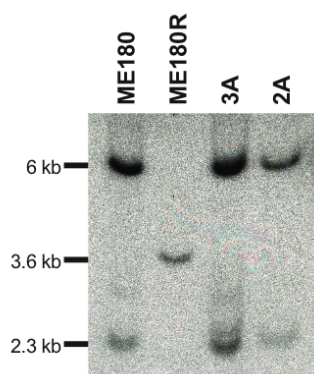


Figure 2.11: Southern hybridization analysis of integrated HPV68b in ME180, ME180R, ME180-2A and ME180-3A. BamHI-digested genomic DNA from ME180, ME180R, 2A and 3A were hybridized with an HPV68b DNA probe (L1-L2-URR-E6-E7). The sizes of the hybridized fragments are indicated.

2.1.4 Cloning and sequencing of a complete HPV68b genome from a CIN2 lesion

Because a complete sequence of HPV68b could not be assembled from the two incomplete copies of HPV68b(int)-ME180 and also was not available from the literature, it was aimed to clone and sequence a complete HPV68b genome from a cervical lesion. Through the DKFZ-CGE collaboration, an HPV68-positive DNA sample of a CIN2 lesion was provided. The DNA was amplified with phi29 DNA polymerase, in the process called multiple displacement amplification (MDA), because the original DNA amount was not enough for further analysis. The original DNA of this CIN2 sample is referred to as original-CIN2, and the amplified product as MDA-CIN2.

To determine whether the CIN2 sample contained HPV68b and not HPV68a, a partial sequence of ORF L1 was determined. Sequence differences between HPV subtypes are in the range of 2-10% in the L1 gene, whereas variants differ only by less than 2% (de Villiers et al, 2004). The 3064-bp region covering the HPV68 ORFs L2-L1 was amplified from MDA-CIN2 DNA (Figure 2.12). The product was cloned and sequenced from one end with the reverse primer as the sequencing primer. The sequence of 798 bp covering part of ORF L1 (Figure 2.12) was used for comparison with the corresponding region of HPV68b(int)-ME180 and HPV68a (DQ080079). The 798-bp sequence of the CIN2 sample showed 99% sequence identity to HPV68b(int)-ME180 3'copy (793 out of 798 nucleotides) and 91% to the HPV68a (732 out of 798 nucleotides). This result demonstrated that the CIN2 sample contains subtype HPV68b. In addition, analysis by the HPV Luminex assay (Schmitt et al, 2006) indicated that the CIN2 DNA contains HPV68b

as single HPV type. Therefore, the CIN2 DNA was used as source for isolation of the complete HPV68b genome.

To amplify the complete HPV68b genome from the CIN2 sample, two primer pairs were initially used (Figure 2.13, panel A) with MDA-CIN2 as template, using the Expand Long Template PCR System. The two primer pairs locate about 600 bp apart. From both primer pairs, products of ~8 kb were obtained (data not shown). Attempts to directly clone these two fragments into a TA vector failed. However, it was discovered during the cloning attempts that an EcoRI restriction site is present in the HPV68b-CIN2 genome (Figure 2.13, panel A). After EcoRI cleavage of the ~8 kb PCR products, the EcoRI site could be located in the L1 gene. An 851-bp area covering the suspected EcoRI site was amplified, cloned and sequenced. The sequence revealed an EcoRI site 390 bp downstream of the L1 start codon.

New primers flanking the EcoRI site were designed based on the 851-bp sequence, and were used to amplify the complete HPV68b-CIN2 genome from MDA-CIN2 DNA (Figure 2.13, panel B). The expected ~8 kb product was obtained, but also a smaller product of ~7 kb (Figure 2.13, panel B). The 8 kb and 7 kb fragments were digested with EcoRI for cloning into a pBluescript-KS-Plus (pBS) vector at the EcoRI site. The first attempt failed, probably due to the low concentration of the inserts. Therefore, the EcoRI-digested 8 kb and 7 kb fragments were re-ligated to produce circular DNA. These circular DNAs were amplified with phi29 DNA polymerase, digested with EcoRI and successfully cloned into the pBS vector at the EcoRI site.

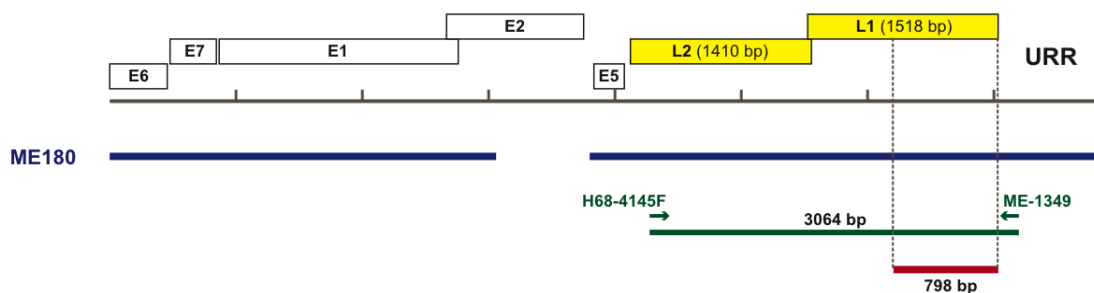


Figure 2.12: Determination of the HPV68 subtype in the CIN2 sample. The HPV68 genome is shown at the top. The genome parts present in ME180 are indicated by the blue lines. The L2-L1 PCR fragment amplified from MDA-CIN2 DNA is shown by the green line. The primers H68-4145F and ME-1349 bind at 1363 bp upstream from the L1 ATG start codon and at 189 bp downstream of the L1 stop codon, respectively. The reverse primer ME-1349 was also used as sequencing primer and a 798-bp partial L1 sequence was obtained (red). The 798-bp sequence of CIN2 was compared with HPV68b-ME180 and HPV68a for subtype determination. The primer sequences are shown in Materials and Methods.

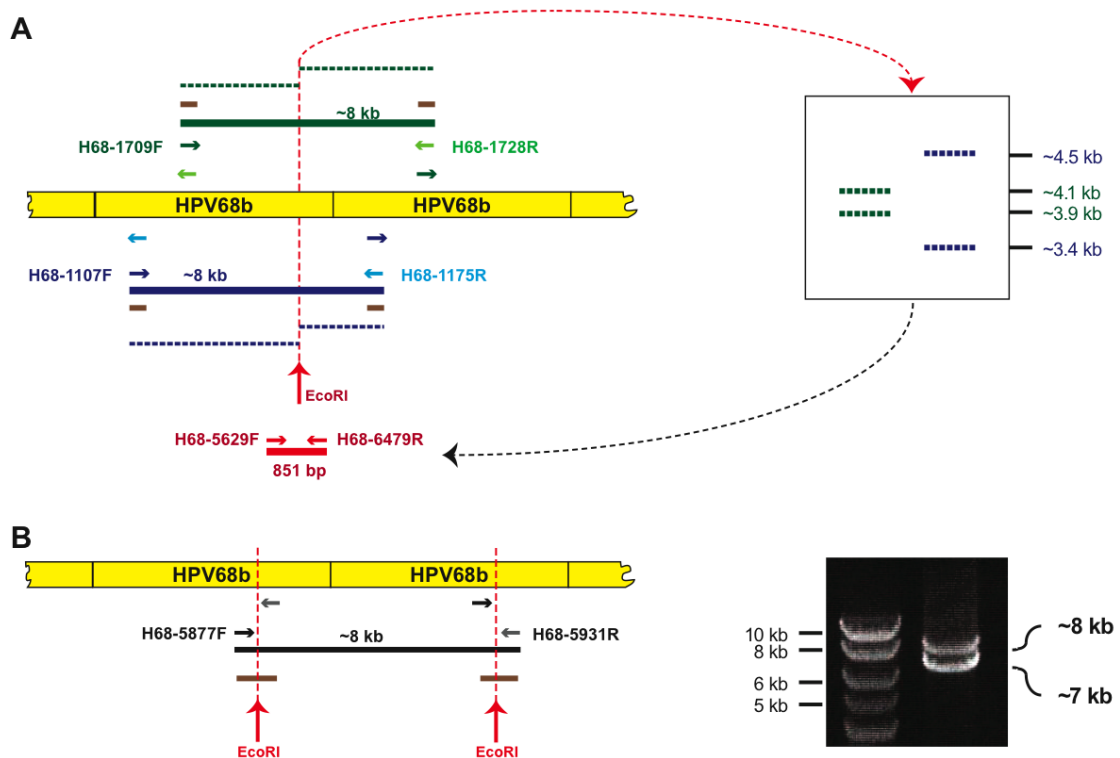


Figure 2.13: Amplification strategies for cloning the complete HPV68b genome from the CIN2 sample. Whole genomic amplification with phi29 DNA polymerase results in concatameric HPV68b DNA (yellow bar) which was used as template molecule for PCR of CIN2 DNA. **Panel A:** Two sets of primer pairs (green and blue arrows) were used for PCR, giving rise to two overlapping full-length genomes (green and blue bars). Primer pair H68-1709F and H68-1728R also produced a 20-bp product, and primer pair H68-1107F and H68-1175R a 69-bp product, indicated by brown bars. The EcoRI site is indicated by the red vertical arrow. On the right, the EcoRI restriction fragments of the two PCR products separated on an agarose gel are schematically shown with their estimated sizes. The 851-bp area (red bar) covering the suspected EcoRI site was amplified with a primer pair binding at 122 bp and 972 bp downstream of the ATG of L1 start codon. **Panel B:** A primer pair (black arrows) flanking the EcoRI site was used for PCR to amplify the full-length genome (black bar). The primers also produce a 54-bp product (brown bars). Agarose gel electrophoresis of the PCR product showed two products of ~8 kb and ~7 kb (right lane). Primer sequences are available in Materials and Methods.

For both 8 kb and 7 kb HPV68b-CIN2 genomes, one positive clone was completely sequenced by Sanger method (A. Hunziker, DKFZ Core Facilities). The nucleotide sequences of the 8 kb and the 7 kb genomes, HPV68b-CIN2 and HPV68b-CIN2-Del, are shown in Appendix A2. The complete HPV68b-CIN2 genome contains 7836 bp. Figure 2.14 shows the HPV68b-CIN2 genome in a circular form, and Table 2.5 the genomic organization. The sequence of the HPV68b-CIN2-Del genome contains 6611 bp and carries a 1229-bp deletion in the E1 gene (Figure 2.15). This 1.2 kb deletion, located between nucleotide positions 1378 and 2608 corresponding to amino acid residues 186 to 595, leads to a frameshift and a premature stop codon at amino acid residue 186. In addition to the 1229-bp deletion, the HPV68b-CIN2-Del genome carries two 1-bp

insertions, one 2-bp insertion, and 11 mutations (Table 2.6). The two single nucleotide insertions in ORFs L1 and L2 would lead to frameshifts. However, these nucleotide differences were not verified by independent cloning/sequencing. The HPV68b-CIN2 and HPV68b-CIN2-Del genomes are illustrated in comparison with HPV68b-ME180 (Figure 2.16).

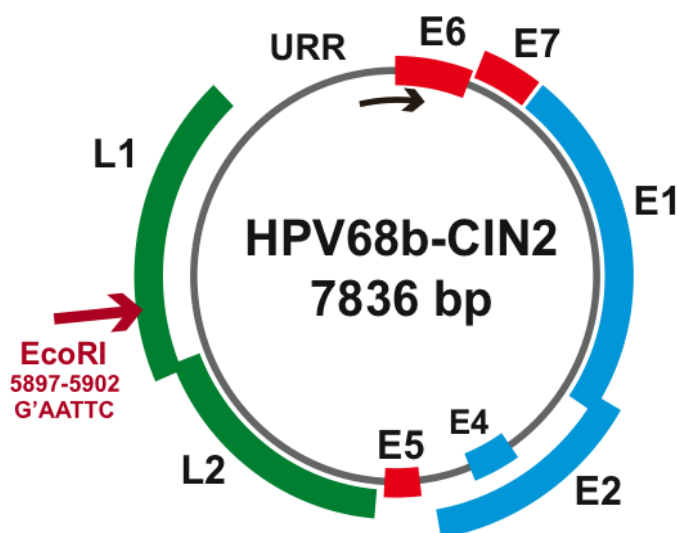


Figure 2.14: The complete HPV68b-CIN2 genome in circular form. The nucleotide positions of each gene and the encoded amino acid residues are given in Table 2.5. The ORFs are shown as red, blue and green bars. The single EcoRI site is indicated by the red arrow. The black arrow indicates the direction of transcription.

Table 2.5: Sequence organization of the complete HPV68b-CIN2 genome.

Gene	ORF* positions	Length (bp)	Protein (aa residues)
E6	1-477	477	158 aa
E7	485-817	333	110 aa
E1	824-2746	1923	640 aa
E2	2673-3785	1113	370 aa
E4	3268-3552 **	285	94 aa
E5	3850-4071	222	73 aa
L2	4118-5527	1410	469 aa
L1	5508-7025	1518	505 aa
URR	7026-7836	811	-

The first nucleotide of the E6 ATG start codon was taken as position 1 for sequence numbering.

* Positions from ATG to stop codon.

** ORF of E4 gene does not start with ATG. It was allocated by comparison with the ORF E4 of HPV68a (DQ080079), with the first position located after the upstream stop codon.

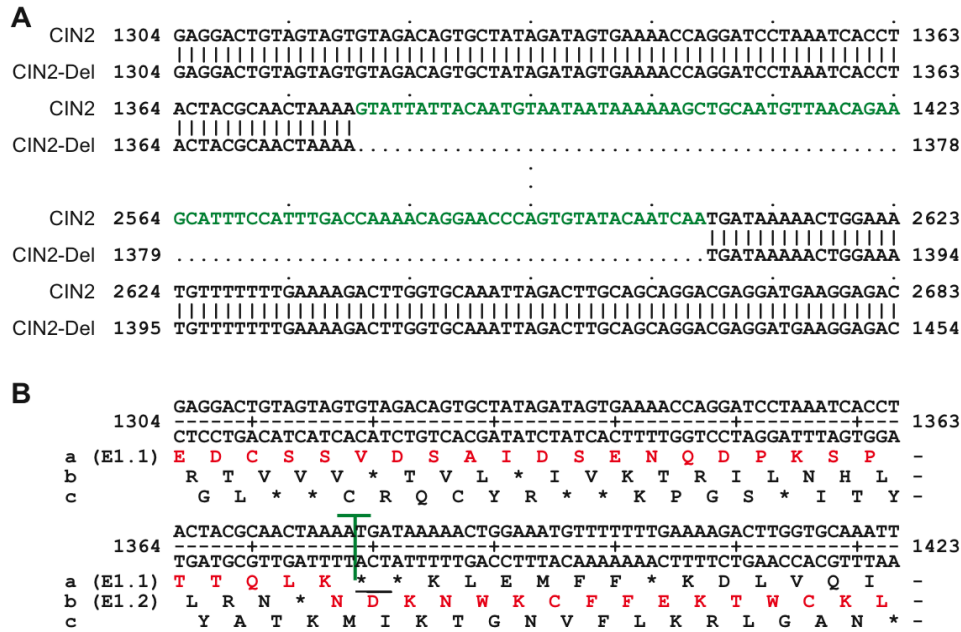


Figure 2.15: Deletion of 1229 bp in ORF E1 of HPV68b-CIN2-Del. Panel A: Nucleotide sequence comparison between HPV68b-CIN2 and HPV68b-CIN2-Del genomes in the left and right deletion boundaries. The nucleotides not present in HPV68b-CIN2-Del are in green. Panel B: Partial nucleotide sequence and translated amino acids of HPV68b-CIN2-Del genome. The 1229-bp deletion (green T-mark) leads to a premature stop codon (underlined) immediately following the deletion site. Red color indicates amino acids translated from normal ORF E1.

Table 2.6: Nucleotide alterations of HPV68b-CIN2-Del in comparison to HPV68b-CIN2.

Positions based on HPV68b-CIN2	Alterations in HPV68b-CIN2-Del	Locations	Effect on proteins
73	T>C	ORF E6	silent mutation
1024	A>G	ORF E1	silent mutation
1379-2607	1229 bp deleted	ORF E1	premature stop codon
3065	T>C	ORF E2	silent mutation
3378	C>T	ORF E2	missense mutation
4670	C>T	ORF L2	missense mutation
5022^23	A-insert	ORF L2	frameshift
5176	C>T	ORF L2	silent mutation
5837^38	G-insert	ORF L1	frameshift
6366	G>A	ORF L1	missense mutation
7307^08	TA-insert	URR	
7486	C>T		
7648	T>C		
7698	T>C		
7778	G>A		

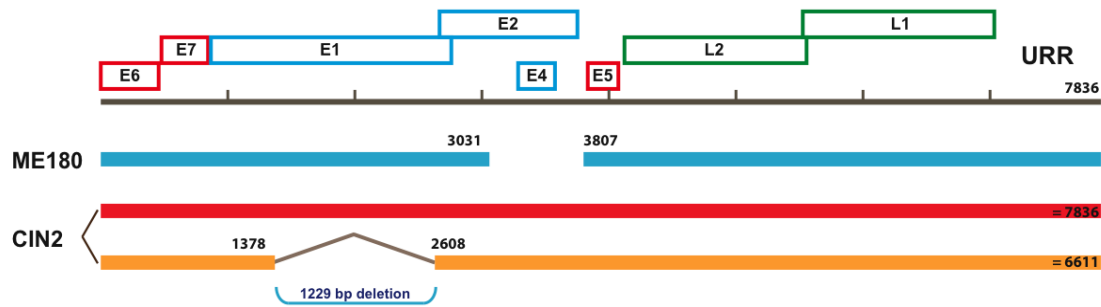


Figure 2.16: HPV68b genomes in ME180 and CIN2. The complete HPV68b genome with the ORFs based on the HPV68b-CIN2 sequence is shown at the top. The genome coverage of the integrated HPV68b in ME180 is shown by the blue lines. The red and orange lines depict the complete and partially deleted HPV68b-CIN2 genomes. The nucleotide positions flanking the deleted regions are given.

Because the HPV68b-CIN2 genomes were obtained through several steps of amplification, it was examined whether the HPV68b-CIN2-Del genome is an artifact or indeed present in the original CIN2 lesion. To clarify this, the area covering the 1.2-kb deletion was amplified from the original CIN2 DNA and from the MDA-CIN2 DNA. A 2.1-kb PCR product was expected from the complete genome, and a 0.9-kb product from the deleted genome. As shown in Figure 2.17, the 0.9-kb PCR product band could be identified in both MDA-CIN2 and original CIN2 DNA. Based on the product amounts of the 2.1-kb and 0.9-kb bands, it seemed possible that the complete HPV68b-CIN2 genome is present in much lower amount in the CIN2 DNA than the partially deleted HPV68b-CIN2-Del genome.

Previously, we had assumed that the HPV68b-CIN2 genome is present as episome in the CIN2 lesion. However, this assumption became disputable after isolation of the HPV68b-CIN2-Del genome. An episomal HPV genome relies on a functional E1 protein for replication and maintenance in the infected cells. Therefore, the question was addressed whether the HPV68b-CIN2-Del genome is present as integrated or episomal DNA. A genomic Southern hybridization was performed using the MDA-CIN2 DNA because of the limitation of the original DNA. EcoRI (single cut) and BamHI (double cut) were used for digestion, before hybridization with the complete HPV68b-CIN2 probe. The hybridization results are shown in Figure 2.18.

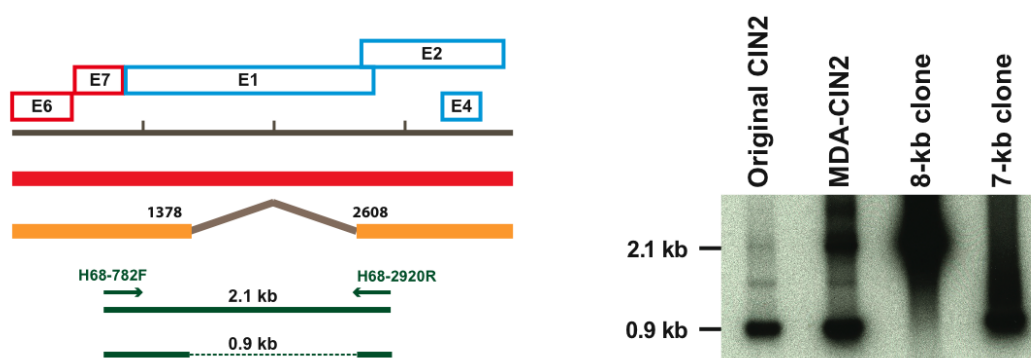


Figure 2.17: Proof of presence of the HPV68b-CIN2-Del genome in the original CIN2 DNA. The PCR covering the deleted region of the HPV68b-CIN2-Del genome is shown on the left. Red and orange lines represent the complete HPV68b-CIN2 and partially deleted HPV68b-CIN2-Del, respectively. The two expected PCR products (green lines) from both HPV68b genomes are indicated with their sizes. The primers are indicated by green arrows. The Southern hybridization of the PCR products from four different templates is shown on the right. Complete HPV68b-CIN2 DNA was used as probe. The 8 kb HPV68b-CIN2 (“8-kb clone”) and 7 kb HPV68b-CIN2-Del (“7-kb clone”) clones were used as positive controls for the 2.1-kb and 0.9-kb PCR products.

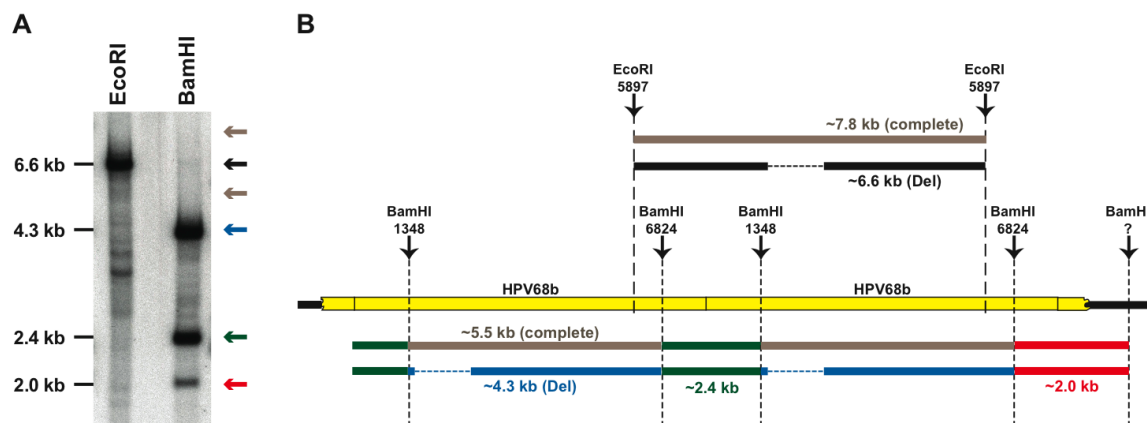


Figure 2.18: Genomic Southern hybridization of MDA-CIN2 DNA. MDA-CIN2 DNA was digested with EcoRI and BamHI, and the filter was hybridized with the complete HPV68b-CIN2 as radiolabelled probe (**panel A**). The ~6.6 kb (black arrow) fragment was expected from EcoRI-digested HPV68b-CIN2-Del genome, and the 4.3 kb (blue arrow) and 2.4 kb (green arrow) fragments from BamHI-digested HPV68b-CIN2-Del genome. The 2.0 kb fragment (red arrow) is assumed to represent a viral-cellular junction fragment. The brown arrows point to positions where the fragments of the complete HPV68b-CIN2 genome would be expected. **Panel B:** Restriction fragments obtained by BamHI and EcoRI digestion of the complete and partially deleted HPV68b genomes. If integrated, the HPV68b genome in CIN2 is assumed to be present in concatameric form (yellow bar). BamHI and EcoRI restriction sites are indicated by vertical black arrows. The color code for the restriction fragments is identical to that used for the horizontal arrows in panel A. Integration junctions are indicated at both ends of the concatameric fragment.

For the complete HPV68b-CIN2 genome, a ~7.8 kb EcoRI fragment, and two BamHI fragments of ~5.5 kb and ~2.4 kb were expected. Neither the 7.8 kb nor the 5.5 kb band were detected after hybridization. For the HPV68b-CIN2-Del genome, a ~6.6 kb EcoRI fragment, and two BamHI fragments of ~4.3 kb and ~2.4 kb were expected. All three bands were observed after hybridization (Figure 2.18). The results indicated that the complete HPV68b-CIN2 genome is present only in minute amount, compared to the HPV68b-CIN2-Del genome. An additional BamHI fragment of ~2.0 kb hybridized to the probe (Figure 2.18, panel A). This fragment could probably represent a viral-cellular junction fragment of integrated HPV68b DNA (Figure 2.18, panel B). Based on this assumption, the HPV68b-CIN2-Del genome is most likely integrated, and thus, the E1 gene is not necessary for genome maintenance. Despite best effort, experimental proof of integration could not be obtained.

Recently, the sequence of a complete HPV68b variant isolated from a patient in China has been published (Wu et al, 2009). Based on the accession number, this HPV68b variant will be named HPV68b-EU918769 in the following. To determine the sequence similarity of HPV68b-EU918769 to HPV68b-ME180 and HPV68b-CIN2, the sequences of these three genomes were compared. The complete genomes of HPV68b-EU918769 and HPV68b-CIN2 have 99% similarity. HPV68b-EU918769 carries a 2-bp deletion in the URR region compared to HPV68b-CIN2. The location of this 2-bp deletion is identical to the 2-bp deletion in the URR of HPV68b(int)-ME180 5'copy (Figure 2.19 and Table 2.7). The nucleotide sequence at this particular location is characterized by different numbers of repeats of the dinucleotide TA (4x, 5x and 6x) among the compared seven HPV68b genomes (Figure 2.19). This feature suggests that this TA-repeat maybe prone to synthesis error by DNA polymerase. Because sequence differences in ORF L1 are used to classify HPV types, subtypes and variants (de Villiers et al, 2004), the ORFs L1 of the three genomes were compared. They have identical length of 1518 bp and 99% sequence identity among each other (Table 2.8). Thus, they are three variants of subtype HPV68b.

Table 2.7: Sizes of genomes, URR and ORFs of HPV68b-ME180, HPV68b-CIN2 and HPV68b-EU918769.

	HPV68b-ME180	HPV68b-CIN2	HPV68b-EU918769
Complete length	7061*	7836	7834
E6 ORF	477	477	477
E7 ORF	333	333	333
E1 ORF	1923**	1923	1923
E2 ORF	incomplete	1113	1113
E5 ORF	222	222	222
L2 ORF	1410	1410	1410
L1 ORF	1518	1518	1518
URR	809 (5'copy) 811 (3'copy)	811	809

* 775 bp deletion in E2 (see Figure 2.16).

** Complete ORF E1 constructed by combination of the HPV68b(int)-ME180 5' and 3'copy sequences.

HPV68b-CIN2	CATGTAATATATATATA..GTTCT	pos. 7291-7312
HPV68b-CIN2-Del1		
	CATGTAATATATATATATAGTTCT	pos. 6064-6087
3'copy of AA13.1 (M73258)	CATGTAATATATATATA..GTTCT	pos. 3400-3421
HPV68b(int)-ME180 3'copy	CATGTAATATATATATA..GTTCT	pos. 10574-10595
HPV68b(int)-ME180 5'copy	CATGTAATATATATA...GTTCT	pos. 3692-3673
HPV68b(int)-ME180R	CATGTAATATATATATA..GTTCT	pos. 3920-3941
HPV68b-EU918769	CATGTAATATATATA...GTTCT	pos. 7291-7310

Figure 2.19: Gap and insertion in URR of HPV68b genomes. Nucleotide sequences of partial URR of seven HPV68b genomes are shown in comparison. The names are indicated on the left and nucleotide positions on the right. The repeated TA is underlined.

Table 2.8: Nucleotide variations in ORF L1 of HPV68b.

Genome	ORF L1 (positions based on HPV68b-CIN2)																No. of nucleotide differences			
	5	6	3	4	5	6	7	8	9	10	11	12	13	14	15	16	HPV68b-CIN2	HPV68b(int)-ME180 5'copy	HPV68b(int)-ME180 3'copy	HPV68b-EU918769
HPV68b-CIN2	T	T	G	C	G	A	G	C	A	C	G	A	A	C	C		-	5	8	9
HPV68b(int)-ME180 5'copy			A				A					G		T	T		5	-	5	8
HPV68b(int)-ME180 3'copy	C	C	A				A				A		G	T	T		8	5	-	11
HPV68b-EU918769			A	T	A	G	A	G	C	A				T			9	8	11	-

2.1.5 HPV68 variant analysis based on URR region

Through the DKFZ-CGE collaboration, eleven HPV68-positive DNA samples from cervical scrapes, including the CIN2 sample, were provided to us. Except HPV68b-CIN2, all other ten samples contain multiple HPV types (Table 2.9). Because the HPV genotyping cannot differentiate between HPV68 subtypes, it was the question whether the DNA samples contained HPV68b or HPV68a. Furthermore, the question was addressed whether the samples belong to already known or new HPV68b variants.

Typically, HPV variants are classified based on sequence differences of the L1 gene (de Villiers et al, 2004). Other regions of the HPV genome have also been analyzed for sequence variations. In case of HPV68, partial URR sequences of forty-one HPV68-positive DNA samples were used to determine the phylogeny of HPV68, and sixteen variants of subtype HPV68b were identified based on the phylogenetic tree (Calleja-Macias et al, 2005).

Because the study of Calleja-Macias and coworkers provides the most extensive source for sequences of HPV68 variants, it was decided to determine the HPV68 sequences in the URR region of the 11 samples (Table 2.9) to compare them with the established 16 variants. The partial URR sequence corresponds to pos. 7279-7769 (491 bp) of HPV68b-CIN2. The sequences of the sixteen HPV68b variants were used as references. One of these variants contains sequences identical to HPV68b(int)-ME180 3'copy. In addition, the recently published complete sequence of an HPV68b variant isolated from China (accession EU918769) (Wu et al, 2009) was also included. For subtype HPV68a, the sequence of accession DQ080079 was used as reference. No further HPV68a sequence has yet been published. To obtain the URR sequences of the 10 additional samples, a region of ~773-bp (pos. 7279-7769 on HPV68-CIN2) covering this segment was amplified from all samples using primers H68-7179F and H68-128R. The PCR products were cloned and sequenced. The sequence of each sample was derived from three clones of identical sequences. The nucleotide sequences from pos. 7279-7769 of all 11 samples in Table 2.9 were compared with the seventeen HPV68b reference variants and the single HPV68a variant using ClustalW. The sequence alignment is shown in Figure 2.20. The nucleotide sequences are shown in Appendix A3.

Table 2.9: HPV68-positive cervical DNA samples.

Sample ID	Sample name	Cytology/ Histology	HPV genotyping*
Reims-06	44667	normal	16, 52, 68 , CP6108
Reims-09	46672	HSIL	16, 54, 68 , 70, IS39
Reims-12	52553	LSIL	39, 51, 52, 61, 68 , 83, CP6108
Reims-14	55178	normal	66, 68
Reims-15	55670	LSIL	6, 16, 31, 35, 40, 53, 61, 68 , (52)
Reims-16	55954	HSIL	16, 58, 68 , (52)
Reims-18	56904	HSIL	16, 55, 56, 68
Reims-28	64767	LSIL	16,51,66, 68 ,69,73,83
Reims-31	66601	ASC-US	16,31,(52),58,59,67, 68 ,72,CP6108
Reims-33	68374	LSIL	16,42,53,61, 68
HPV68b-CIN2	CIN2	CIN2	68b

* HPV genotyping for the ten Reims samples was performed in Reims.

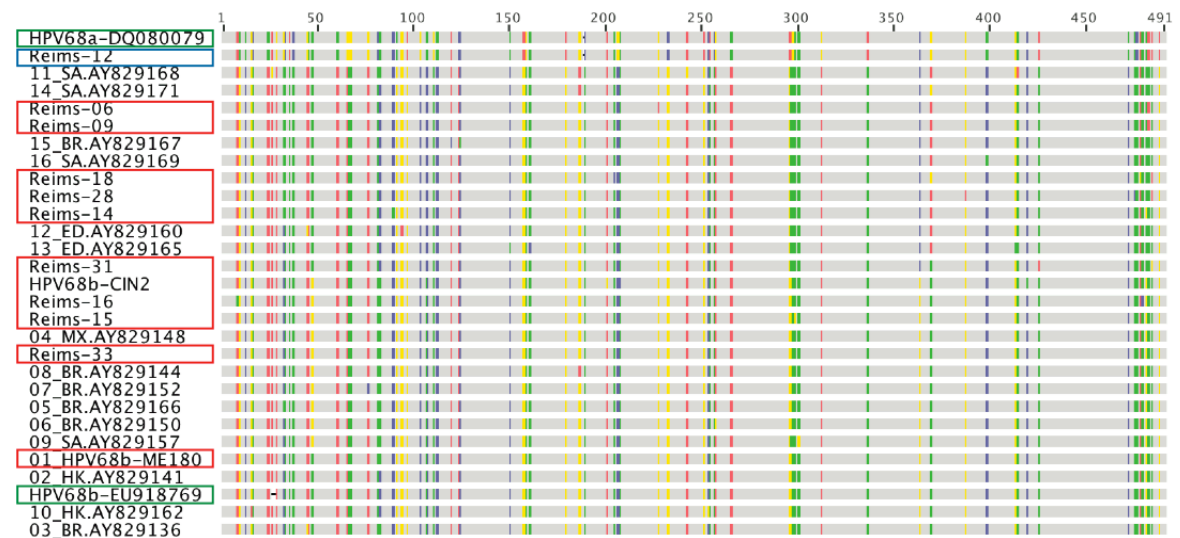


Figure 2.20: Alignment of partial URR sequences of different HPV68 isolates. The 491-bp URR segment (pos. 7279-7769 on HPV68b-CIN2 genome) of eleven clinical samples (Table 2.9) were aligned with the corresponding sequence of HPV68b-CIN2, HPV68b-EU918769 (Wu et al, 2009), HPV68a (accession DQ080079) and sixteen HPV68b variants (Calleja-Macias et al, 2005) using ClustalW. The nucleotide positions are indicated on top. The bases are highlighted at positions of sequence differences. Red represents A, blue C, yellow G, and green T. Gray area represents identical sequences. The sequence names are shown on the left. Sequences with blue or red boxed names were obtained in this PhD work. Except variant 01_HP68b-ME180, the names of all reference HPV68b variants from Calleja-Macias et al start with a number (02-16), followed by the isolate origin (SA: South Africa, BR: Brazil, ED: Edinburg, MX: Mexico, HK: Hong Kong) and the accession number.

The 491-bp sequence of sample Reims-12 is identical to HPV68a-DQ080079, therefore it contains HPV68a. The other ten samples are HPV68b (Figure 2.20). Altogether, 44 nucleotide positions show sequence differences among the analyzed HPV68b samples (Table 2.10). Sample Reims-33 is identical to HPV68b 04_MX.AY829148. The other

samples seem to be new variants due to additional nucleotide polymorphisms compared to the most similar variants (Table 2.10).

To create a phylogenetic tree containing the new sequences with the 16 reference HPV68b variants, the Unweighted Pair Group Method with Arithmetic Mean Algorithm (UPGMA) algorithm was used. This algorithm was also applied in the earlier work (Calleja-Macias et al, 2005). The generated tree is shown in Figure 2.21. In agreement with the previous data (Calleja-Macias et al, 2005), the phylogenetic tree carries two deep dichotomic branches, representing HPV68a and HPV68b. In the HPV68b branch, the variants form two main clusters. The Chinese isolate HPV68b-EU918769 was grouped with the Hong Kong isolate 02_HK.AY829141, thus showing geographic correlation. Based on this tree, one known HPV68a variant, one known HPV68b variant and nine new HPV68b variants were identified among the eleven HPV68-positive samples from France.

Table 2.10: Nucleotide polymorphisms in the URR of HPV68b variants

[illegible]

* Based on Figure 2.20 (first row) and the sequence of the complete HPV68b-CIN2 genome (second row).

** Obtained from (Calleja-Macias et al, 2005).

“-“ indicates a gap.

Green highlight indicates variations not present in variants 01-16.

Green highlight indicates variants not present in variants of Ref. Variants assignment of the ten DNA samples (bottom part) is based on sequence similarity to one of the reference variants, where “new” indicates a new variant and the bracketed numbers refers to the most similar reference variant.

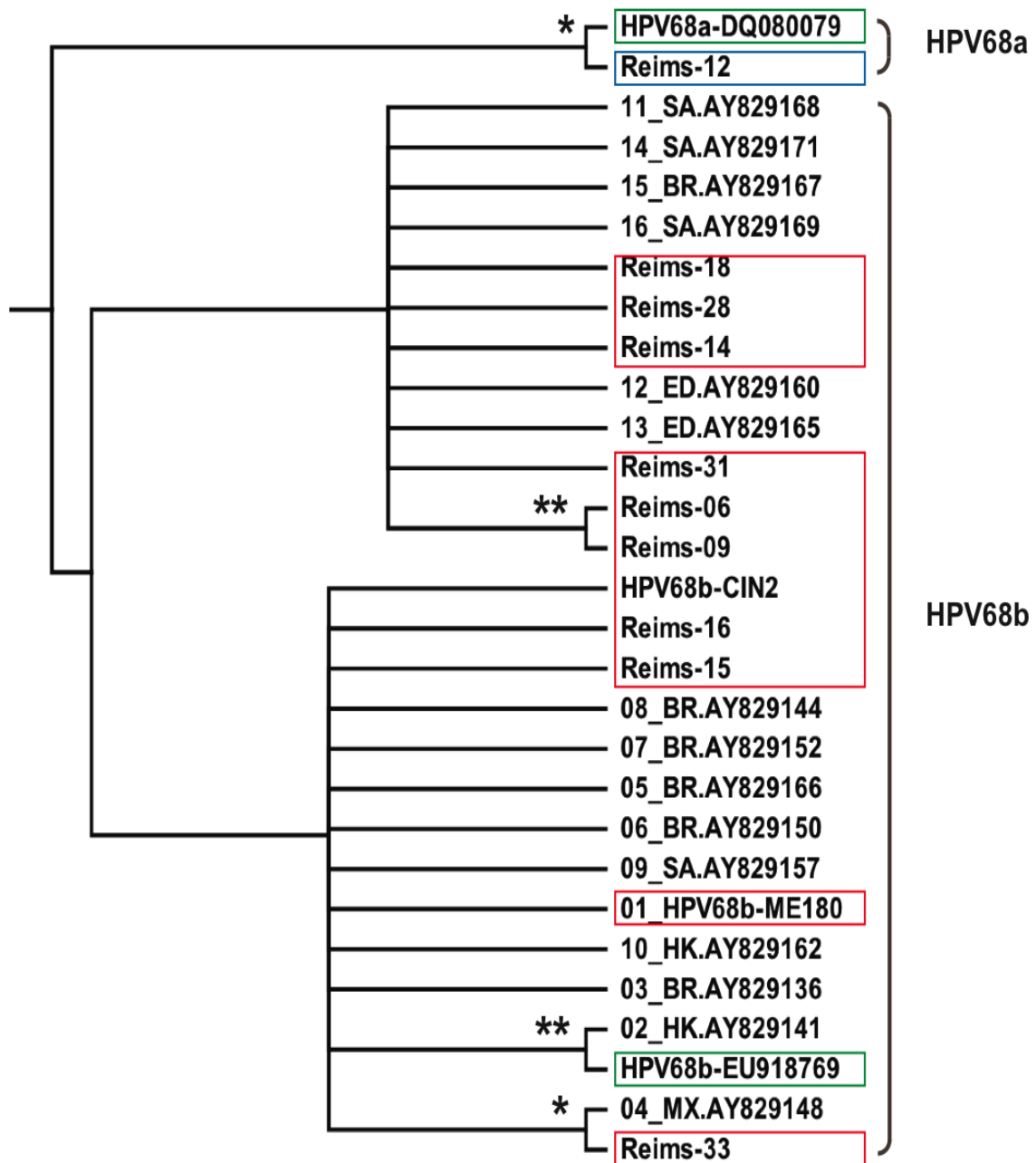


Figure 2.21: Phylogenetic tree created from partial HPV68 URR sequences. A rooted phylogenetic tree was created from the sequence alignment of HPV68 URR sequences (pos. 7279 to 7769 based on HPV68b-CIN2 genome), using UPGMA algorithm. The sequences obtained in this PhD work are bracketed in red and blue, sequences from the database in green, whereas the sixteen previously published HPV68b variants (Calleja-Macias et al, 2005) are shown without brackets. * indicates identical sequences. ** indicates sequences are not identical.

2.2 Analysis of HPV16 sequences generated by the ASP16 strategy

The ASP16 strategy was developed by Bo Xu as a novel strategy for determination of HPV16 integration sites in cervical lesions (Xu, 2010). The ASP16 strategy employs one of the next generation sequencing technologies, Roche/454 GS-FLX pyrosequencing (<http://www.454.com/>). With the Roche/454 GS-FLX format used, each ASP16 sequencing run can generate around 200,000 sequence reads. For this enormous amount of sequence data, it is necessary to apply computer power to assist in sequence data analysis. With my knowledge in computer programming, I collaborated with Bo Xu in the ASP16 project, and developed a set of computer programs for automatically analyzing sequence data generated by Roche/454 GS-FLX sequencing. The development of these programs is described in section 2.2.1. After Bo Xu had finished his PhD work with two ASP16 experiments completed, I continued to optimize the ASP16 strategy and conducted another two ASP16 experiments (ASP16-3 and ASP16-4). The results of these two experiments are described in section 2.2.2.

2.2.1 Development of computer programs for ASP16 analysis

2.2.1.1 Designing computer program strategies and platforms

In ASP16 strategy, the DNAs from cervical scrapes are amplified by GenomePlex whole genome amplification. The HPV16-containing DNAs are enriched by linear amplification using a set of HPV16-specific primers spanning the HPV16 E1-E2 region to produce single-stranded fragments, which are then amplified in multiplex PCR using bipartite forward primers and a tripartite reverse primer (Figure 2.22). There were 24 different 4-nt barcodes designed for the ASP16 experiments. Each barcode is unique for an individual DNA sample.

The primary aim of the ASP16 strategy is the identification of viral-cellular junction sequences. Considering the enormous number of sequence reads in each ASP16 experiment, however, it is not practical to perform human database blasting for all sequence reads. It was decided that the program sets would not perform this task, but narrow down the sequence reads with possible integration junctions by classifying the sequences into small groups using different criteria.

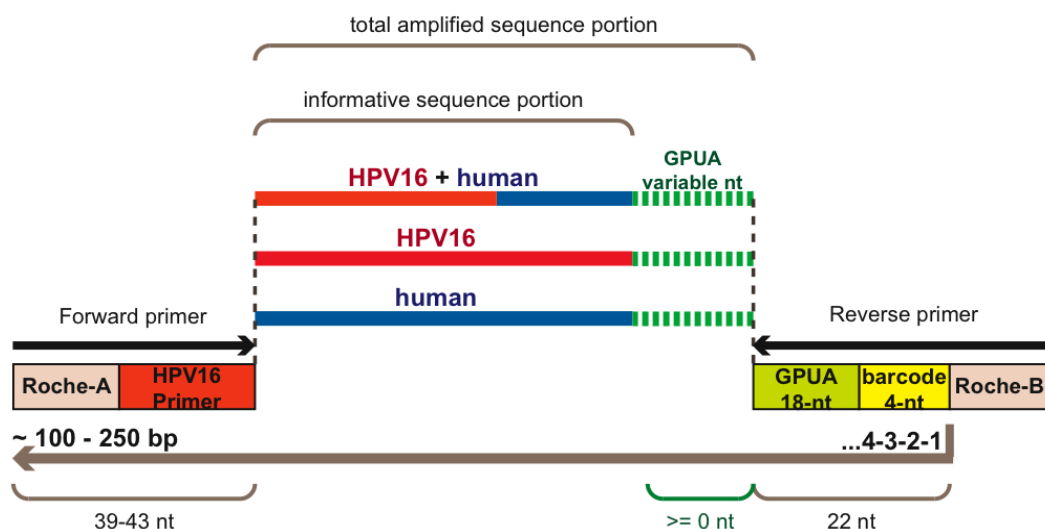


Figure 2.22: Principal structure of amplicon fragments for Roche/454 pyrosequencing in ASP16 experiments. The amplicon fragments are amplified using primer pairs (black arrows) whose compositions are indicated. Roche-B is used as the sequencing primer. The individual sequence read starts with the 4-nt barcode, followed by the minimum length of 18-nt GenomePlex universal adapter sequence (GUA), then by variable amplicon sequence length (pure HPV16, pure human DNA, or HPV16 with human DNA), and ends with RA_HP16 primer. The arrow at the bottom indicates the direction of sequencing. The dashed green line represents a variable length of variable GUA sequence, ranging from 0 up to ~100 nt (see text).

The computer programs were designed to be able to

- (1) sort sequence reads into sample groups where they originated from, according to the 4-nt barcode,
- (2) eliminate sequence reads which do not contain HPV16 DNA sequence (pure human DNA),
- (3) classify sequence reads into groups based on their HPV16 primers,
- (3) generate alignments of grouped sequence reads,
- (4) perform basic statistical analysis, and
- (5) generate outputs in suitable file formats.

The program strategy is summarized in Figure 2.23. The input data is a text file, containing FASTA format sequence reads generated by the Roche/454 GS-FLX pyrosequencing. The program strategy starts by reading the sequence data from the input file. The sequence reads are analyzed one by one. Eventually, alignments of appropriately grouped sequence reads are generated and they may be viewed and edited with any sequence editor software. These alignments are helpful for sequence data visualization in order to search for integration junctions. The basic statistical outputs are generated as tab-delimited files and can be opened in Microsoft Excel or similar software.

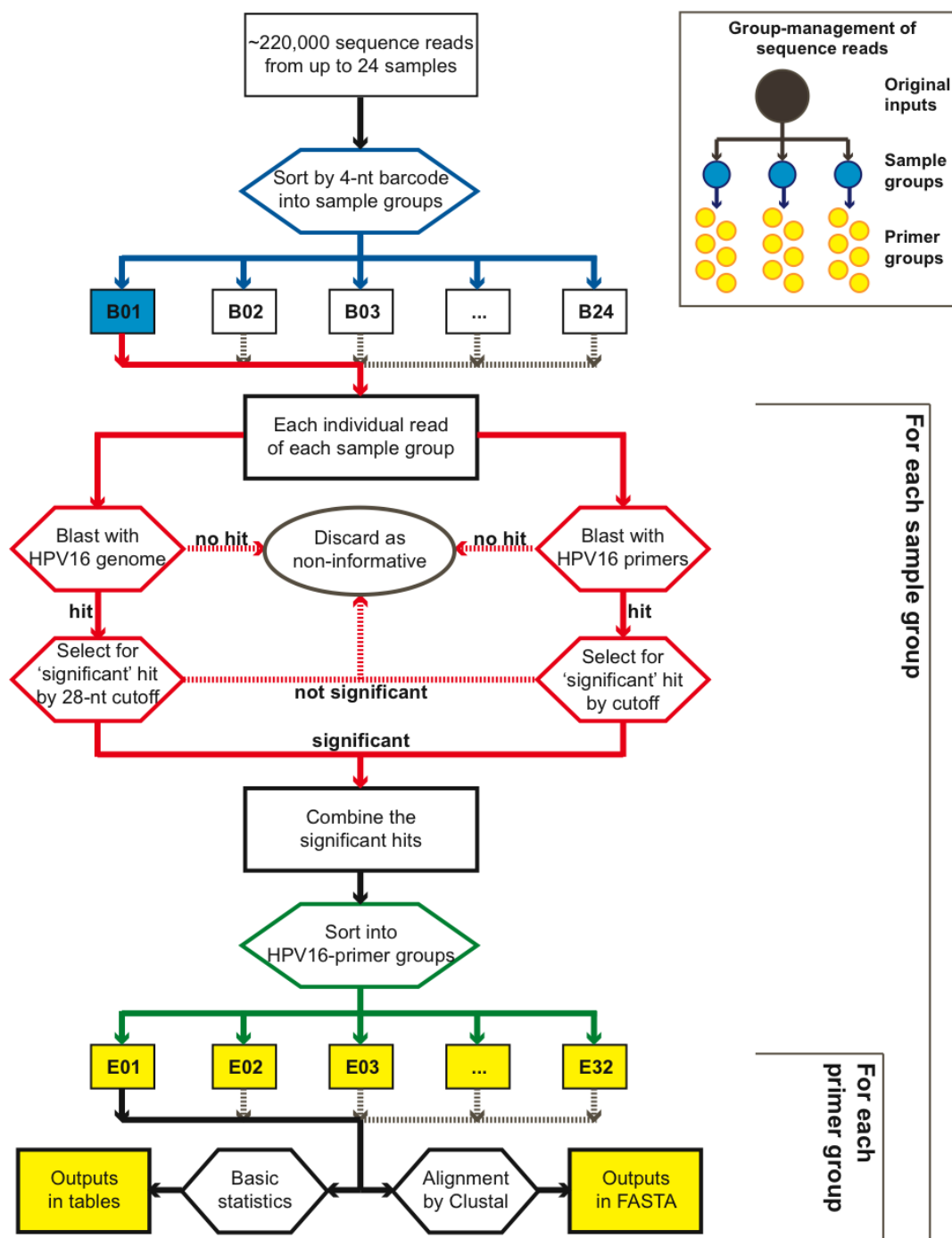


Figure 2.23: Program strategy for ASP16 sequence analysis. From the overall input of up to 220,000 sequence reads from Roche/454 GS-FLX, each sequence read is sorted into one of 24 sample groups according to the 4-nt barcodes (B01, ..., B24). Each barcode is unique for an individual DNA sample. For each sample group, each sequence read is blasted with HPV16 reference genome and a set of complete Roche-A_HPV16 (RA_HPV16) primers. If there is no hit, the sequence is discarded. Otherwise, the hit sequence is filtered with appropriate cutoffs (see text). If the sequence meets the cutoff conditions, it is regarded as significant and processed further, otherwise discarded as non-informative. Significant sequences determined by both HPV16 genome and HPV16 primers are combined, and sorted into HPV16 primer groups (E01, ..., E32). For each primer group, the sequence reads are aligned with ClustalW and the alignments are written as FASTA outputs. Basic statistics are also calculated and outputs are generated in table-format files.

This designed program strategy does not include automatic identification of human DNA in the sequence reads. The main reason originates from the complication of GPUAs sequences. The GUA sequence was incorporated into the reverse primer used to produce the amplicon and has the minimum length of 18 nt. However, as the sequence data of ASP16-1 and ASP16-2 showed, the GUA can be longer with GT-rich unpredictable sequence variations of up to ~100 nt (Xu, 2010). Because of this unpredictability, it was not always possible to define whether the GT-rich variable sequence belongs to GUA or human sequence (integration junction) without blasting it to the human sequence database. Blasting of such GT-rich sequence to the human database usually takes at least a few minutes and does not always come up with a satisfactory result. Therefore, this procedure is not integrated in the analysis process.

With the main computer program strategies decided, the next task was to decide on a suitable computer platform. Perl is a high-level computer programming language favored in the bioinformatics community (<http://www.perl.org>). For this study, it was chosen as the computer programming language because of its efficiency, simplicity, rapid prototyping, and ease-of-use for object-oriented programming. All computer programs were executed under Perl application version 5.8.8, and primarily designed to run under Mac operating system (Mac OS). With some modifications, they may also run under other Unix-based environments, such as Linux. The programs are not suitable for running under Windows operating system because of memory- and resource-handling problems, which will result in substantially slower program process.

2.2.1.2 Computer program algorithms and their tasks

Four sets of programs were written. Key algorithms of these program sets are described in the following. It is intended to show how the programs proceed and how they make decisions. The program source codes line-by-line will be explained in the Appendix A8.

Program set 1: Sort sequences into sample groups by barcodes

The first program set consists of a single program that sorts total sequence reads by their barcode sequences into sample groups (Figure 2.24). It consists of a single program named `PROG_SET_1_mac.pl`. A single text file, containing complete

sequence reads from the Roche/454 GS-FLX sequencing in FASTA format, is given as an input. The outputs of program set 1 are 24 groups of sequence reads, each group corresponding to an individual DNA sample.

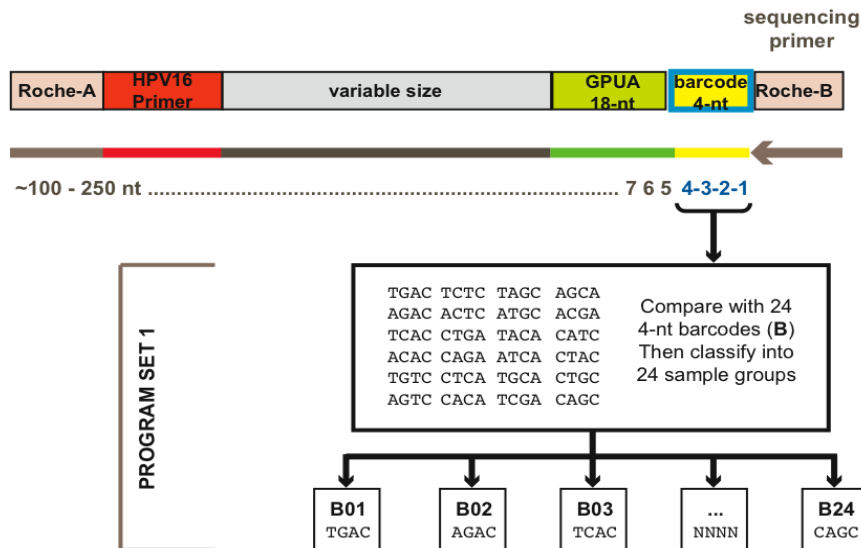


Figure 2.24: Scheme of program set 1. The principal structure of amplicon fragments for the Roche/454 GS-FLX sequencing is shown at the top. Roche-B is used as the sequencing primer. The first four nucleotides of the sequence read correspond to the 4-nt barcode. Each barcode represents an individual DNA sample. This program set compares the first 4 nucleotides of each sequence read with a collection of the twenty-four 4-nt barcodes used in ASP16. The sequences are sorted into 24 barcode groups, based on the matching 4-nt sequences.

Program set 2: Selection of informative sequences

The second set selects informative sequence reads by blasting and filtering steps, sorts the sequence reads into individual HPV16 primer groups, calculates basic statistical data, and delivers the outputs. It consists of nine programs, with the following names:

```
set2_p2mac.pl
set2_p3mac.pl
set2_p4mac.pl
set2_p5mac_a_noCutoff.pl
set2_p5mac_d_28bpCutoff.pl
set2_p5mac_edit.pl
set2_p6mac_before.pl
set2_p6mac_fasta.pl
set2_p7mac.pl.
```

The input to this program set is a group of sequence reads of an individual sample group (also called barcode group) from the outputs of program set 1. Set 2 will select informative sequence reads and group them into HPV16 primer groups. The tasks are divided into three steps: (1) blast with the complete HPV16 reference genome (HPV16R) database and filter with cutoffs, (2) blast with HPV16 primer database, sort into HPV16 primer groups and filter with cutoffs, and (3) combine results of the two database blasts, analyze possible HPV16 breakpoints and report data in Tab-delimited format. The scheme of program set 2 is shown in Figure 2.25.

Set 2 – step 1: Blasting with HPV16R, parsing information and filtering

With regard to the HPV16 genes E1 and E2, the pyrosequencing reads contain the sequence information of the minus (or antisense) strand. For conversion into plus (or sense) strand, the program set generates reverse-complementary sequences of all input sequences (Figure 2.25). This produces sequences starting from Roche-A sequence and ending with the 4-nt barcode sequence. Each reverse-complemented individual sequence is blasted with HPV16R database (`hvp16`), using `blastall` program with the following defined values: `program=blastn`, `word-size=11`, `DUST-filter=off`. Each blast result is saved as a text file. There are three possibilities of the blast results: (1) no hit, (2) one hit, (3) more than one hit. “Hit” refers to a match of a query (ASP16 sequence read) to a region on a sequence in the database. More-than-one-hit indicates matches on more than one location on the HPV16R genome.

After all the sequences have been blasted, the program reads information from each blast result file. Sequences with no hit are discarded. The program parses information from the blast result files of the sequences that have one or more hit. Figure 2.26 shows an example for information parsing.

For sequences with more than one hit, the programs check whether the second hit contains useful information for further analysis. The data of the second hit are considered useful and are collected only if the matched region is at least 94% similar (`%_match`) and covers at least 15 nt (`scoreBP`). Otherwise, the second hit is considered non-informative and only data from the first hit are collected.

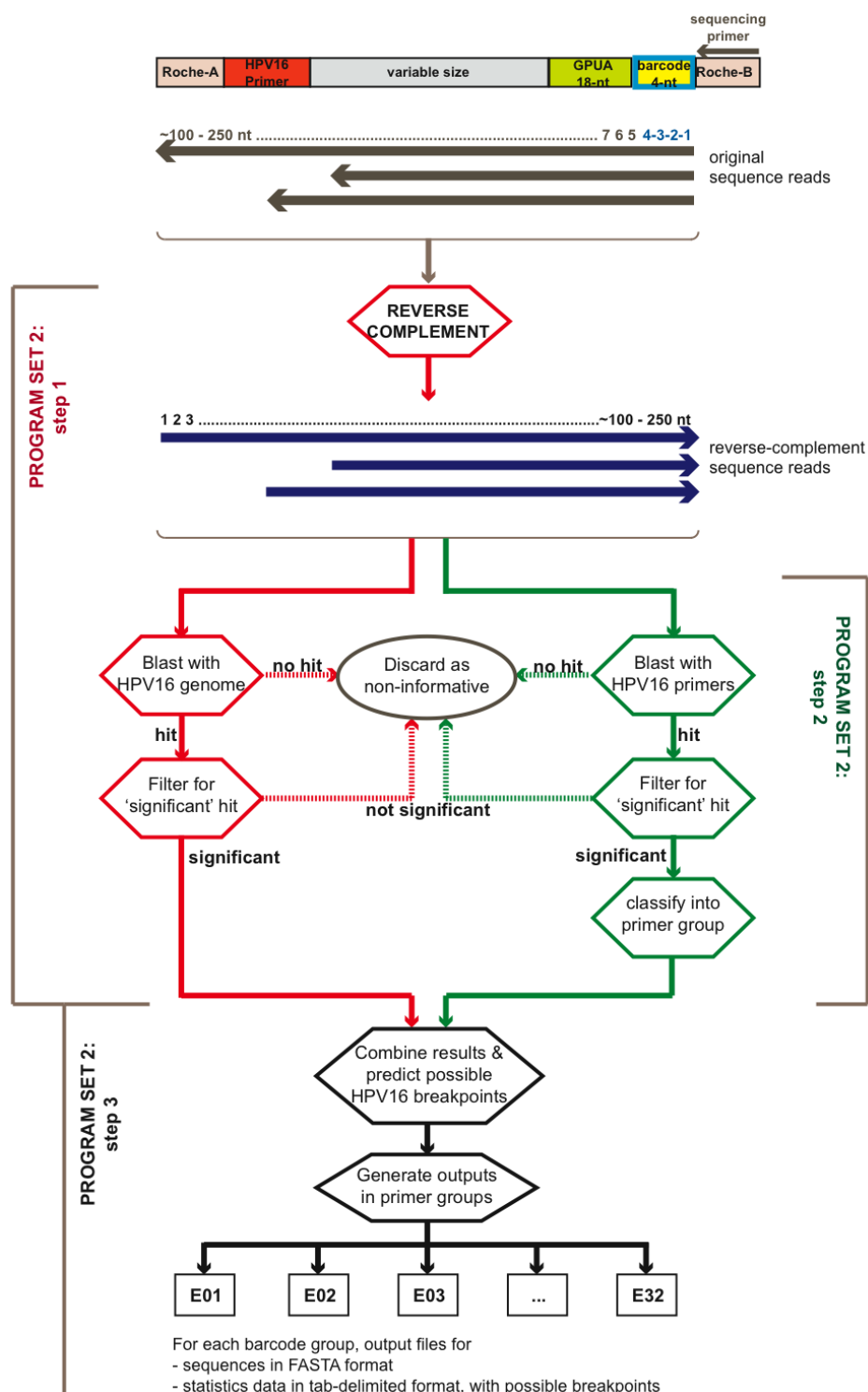


Figure 2.25: Scheme of program set 2. Program set 2 is divided into three steps. In step 1 (left side), the reverse complementary sequences of the original sequence reads are generated and used for further analysis. The reverse-complementary sequences are blasted with HPV16R genome, and sequences with significant HPV16R matches are selected through filtering processes. In step 2 (right side), the reverse-complementary sequences are blasted with RA_HPV16 primer database. The sequences with significant primer matches (see main text for details) are selected through filtering processes, and classified into primer groups. In step 3, significant sequence lists from step 1 and 2 are combined. For each primer group, the possible HPV16 breakpoints are predicted, and outputs are written in tab-delimited format. A file containing the complete sequences of each primer group is generated in FASTA format.

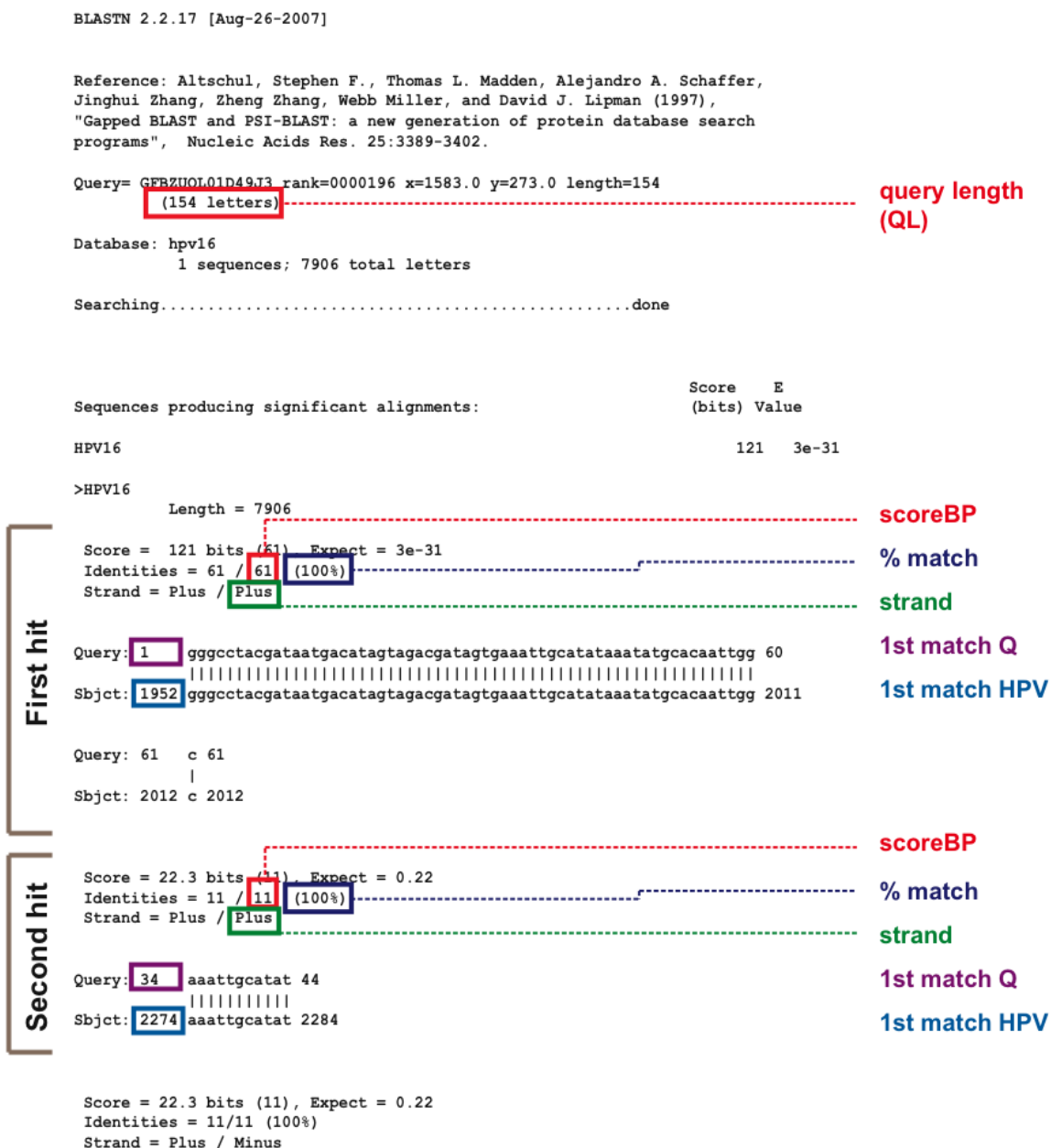


Figure 2.26: Parsing of information from HPV16R blast result. An example of a HPV16R blast result with more than one hit is shown. Query refers to the ASP16 sequence read and Sbjct is HPV16R. For HPV16R blast with one or more hit, six data are collected as marked in colored boxes from the first two hits. *Query length* (QL) is the length of the ASP16 sequence read. *ScoreBP* is the length (nt) of the area where query matches to HPV16R. *% match* is the percentage of similarity between the query and HPV16R. *Strand* refers to the HPV16R matching strand. *1st match Q* is the position on the query strand where matching to HPV16R starts. *1st match HPV* is the position on HPV16R matching to 1st match_Q position. ScoreBP of the first hit is used in the filtering step for significant HPV16R match decision.

The informative sequences should contain a significant length of HPV16R sequence to ascertain that the sequence read is not an unspecific amplicon. This is achieved by the filtering process. A cutoff value of 28 nt was defined for the filtering step. This means that the sequence read is considered significant if the *scoreBP* value of the first HPV16R hit is at least 28 nt. The list of significant sequences are collected and stored for further analysis in the third step.

Set 2 – step 2: Blasting with HPV16 primers, sorting to primer groups and filtering

Each reverse-complementary sequence, generated in step 1, is blasted with the RA_HP16 primer database (EEpyro3), using `blastall` program with the following defined values: `program=blastn`, `word-size=11`, `DUST-filter=off`. Each blast result is saved as a text file. There are three possibilities of the blast results: (1) no hit, (2) one hit, (3) more than one hit. More-than-one-hit indicates matches to more than one HPV16 primer sequences.

Sequences with no hit to RA_HP16 primer are discarded. For sequences with hit(s), there are six possibilities of how the query can match to a RA_HP16 primer. These are outlined in Figure 2.27.

The program set determines, based on these six cases, whether each of the sequence reads has a significant primer hit. The parsed values from RA_HP16 primer blast results are used to evaluate the conditions of the hit. The decisions are made according to different criteria that are listed in Table 2.11. There are three categories of significant RA_HP16 primer matches. Category 1 includes matches as in case A1. Category 2 includes matches as in case A2. Category 3 includes matches as in cases B1, B2, C1 and C2. A sequence read is defined as having a significant RA_HP16 primer match if it meets all the required conditions in one of these three categories. Otherwise, it is declared as having no match to any RA_HP16 primer and excluded from further analysis. The sequences with significant RA_HP16 primer hits are sorted into groups according to the RA_HP16 primers, and stored for the next analyzing step.

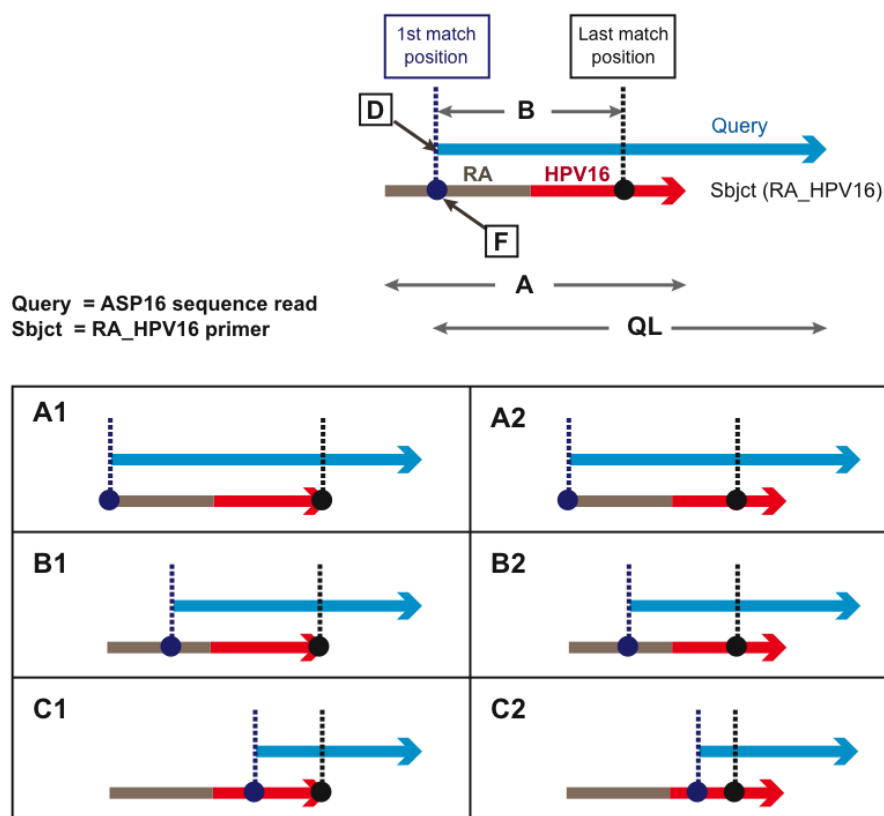


Figure 2.27: Matching possibilities between ASP16 sequence read and a RA_HP16 primer. There are six possible matching patterns between an ASP16 sequence read and a RA_HP16 primer. Cases A1 and A2 show ASP16 reads matching from the beginning of the primer. Cases B1 and B2 show ASP16 reads missing the 5' end of the primer. Cases C1 and C2 show ASP16 reads matching only in the HPV16 part of the primer. In cases A1, B1 and C1, ASP16 reads match until the 3' end of the primer. In cases A2, B2 and C2 the match terminates within the HPV16 part of the primer. QL: query length(nt). A: Sbjct length (nt). B: matching area (nt). D: first matched position on the query. F: first matched position on Sbjct. C: percent similarity (not illustrated). E: strand of Sbjct which the query matches to (not illustrated).

Table 2.11: Criteria for the selection of significant RA_HP16 primer match

Criteria	Criteria*	significant RA_HP16 primer match		
		Category 1 (A1)	Category 2 (A2)	Category 3 (B1,B2,C1,C2)
The matching strands of the query and the primer sequence are in the same orientation	E = Plus	+	+	+
The matching similarity is at least 90%	C >= 90%	+	+	+
The query matches the complete primer length	B >= A	+		
The query does not match the complete primer length	B < A		+	+
The query matches from the 5' end of the primer	F = 1	+	+	
The query does not match from the 5' end of the primer	F != 1			+
The query does not include the Roche-A part of the primer	(B+F) > 20			+
The query does not match the 3' end of the primer, max. 4 nt	(A-(B+F)) > 5		+	+

* The criteria abbreviations are explained and illustrated in Figure 2.27.

Set 2 – step 3: combining blast results, predicting breakpoints and exporting the outputs

In this step, the lists of sequences with significant HPV16R matches and significant RA_HP16 primer matches are combined and compared (see Figure 2.25). Only sequences that are present in both lists are processed further. Otherwise, they are declared as non-informative and discarded. At the same time, the program counts the numbers of sequence reads in each primer groups for both significant-only sequence reads (filtered by 28-nt cutoff in step 1) and total sequence reads (without the 28-nt cutoff filtering). The resulting sequence numbers are exported in two tab-delimited format files, one belonging to no-cutoff and the other belonging to 28-nt cutoff.

In this step, the program estimates the longest position on the HPV16R genome that each sequence read is matched to. This position is informative for hinting at a possible HPV16 breakpoint, and called “last-HPV16-match position”. Using parsed information from HPV16R blast results in step 1 (see Figure 2.26), the values of the first and second (if exists) HPV16R hits of each sequence read are analyzed.

For the sequence reads with a single HPV16R hit, the data of this single hit are used for estimation of the last-HPV16-match position. For sequence reads with more than one HPV16R hits, the program evaluates whether to use the information from both hits or only one of the hits. The decisions are made based on positions and lengths of matching HPV16 area of each hit. All categories for sequences with two hits are explained in Figure 2.28.

After all sequence reads are sorted into appropriate categories, the program reports the selected data in tab-delimited format files for each sample group. Figure 2.29 shows an example of these data, opened in MS Excel. The data are present in 13 columns as explained in Table 2.12. For each DNA sample, a possible HPV16 breakpoint may be predicted by looking at the data in column F. If HPV16 integration exists, a number of sequence reads should have an identical HPV16 position in column F.

















Category	illustration	Explanation and usage
1		The first hit ends before the second hit.
1.1		The first hit starts at the same HPV16 position as the second hit. The data of the second hit are used.
1.2a		The first hit starts earlier than the second hit. Both hits overlap. The second hit is ≥ 15 nt and has $\geq 94\%$ similarity to HPV16R. The data of both hits are combined as a single match.
1.2b		The first hit starts earlier than the second hit. Both hits overlap. The second hit is < 15 nt or has $< 94\%$ similarity to HPV16R. The data of the first hit are used.
1.2c		The first hit starts earlier than the second hit. Both hits do not overlap. The second hit is ≥ 15 nt and has $\geq 94\%$ similarity to HPV16R. The data of the both hit are combined as a single match, with or without a gap.
1.2d		The first hit starts earlier than the second hit. Both hits do not overlap. The second hit is < 15 nt or has $< 94\%$ similarity to HPV16R. The data of the first hit are used.
1.3		The first hit starts after the second hit. The data of the second hit are used.
2		The first hit ends at the same position as the second hit.
2.1		The first hit starts at the same position as the second hit. The data of the first hit are used.
2.2		The first hit starts earlier than the second hit. The data of the first hit are used.
2.3a		The first hit starts after the second hit. The second hit is < 15 nt or has $< 94\%$ similarity to HPV16R. The data of the first hit are used.
2.3b		The first hit starts after the second hit. The second hit is ≥ 15 nt and has $\geq 94\%$ similarity to HPV16R. The data of the second hit are used.
3		The first hit ends after the second hit.
3.1		The first hit starts at the same position as the second hit. The data of the first hit are used.
3.2		The first hit starts earlier than the second hit. The data of the first hit are used.
3.3a		The first hit and second hit overlap. The second hit is ≥ 15 nt and has $\geq 94\%$ similarity to HPV16R. The data of both hits are combined as a single match.
3.3b		The first hit and second hit do not overlap. The second hit is ≥ 15 nt and has $\geq 94\%$ similarity to HPV16R. The data of both hits are combined as a single match, with or without a gap.
3.3c		The first hit and second hit overlap. The second hit is < 15 nt or has $< 94\%$ similarity to HPV16R. The data of the first hit are used.
3.3d		The first hit and second hit do not overlap. The second hit is < 15 nt or has $< 94\%$ similarity to HPV16R. The data of the first hit are used.

Figure 2.28: Different cases of double HPV16R blast hits of a ASP16 sequence read. This figure shows different cases of two HPV16R hits from a ASP16 sequence read. Each line in the illustration column represents a matched area on the HPV16R genome. For sequences with only one HPV16R hit (not shown in this figure), the data of this hit are used for breakpoint analysis. For sequences with two hits, the program decides, based on 16 different cases, which data to be used. Blue represents the first hit, and green the second hit. Continuous lines indicate the data of the hit are used for breakpoint analysis. Dashed lines indicate the data of that hit are not used for the analysis. The lengths and positions of the blue and green lines imitate the relative HPV16 matching length and actual positions on HPV16 genome respectively. Each case is explained in the last column. The usage of hit data is indicated.

A	B	C	D	E	F	G	H	I	J	K	L	M
barcode	primer	SeqNo	Seq Length	Last hit Position Query	Last hit Position HPV	Match nt	Last 2nd-hit Position Query	Last 2nd-hit Position HPV	Match nt 2ndHit	Overlap 1st & 2nd hits	Original Seq Name	Reverse complement sequence
1	E02	38	67	36	1312	36	none	none	none	none	>GFBZUQL01	AAGTGGAAAC
1	E02	299	61	29	1303	29	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	328	77	58	1332	58	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	331	68	47	1322	41	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	439	60	28	1302	28	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	459	75	37	1311	37	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	619	72	36	1310	36	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	644	91	42	1316	42	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	676	72	42	1316	42	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	737	70	38	1312	38	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	807	68	36	1310	36	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	863	64	34	1308	34	none	none	none	none	>GFBZUQL01	TGAAGTGGAA
1	E02	924	69	42	1316	42	none	none	none	none	>GFBZUQL01	TGAAGTGGAA

Figure 2.29: Example of a tab-delimited output from program set 2. For each DNA sample group, a tab-delimited file is exported by program set 2. The file is opened in MS Excel. For further explanation, see Table 2.12.

Table 2.12: Explanation of tab-delimited output from program set 2 as shown in Figure 2.29.

Column	Title	Description
A	barcode	barcode number
B	primer	RA_HP16 primer group, assigned to the sequence read
C	SeqNo	sequence ID, assigned by this program set
D	SeqLength	sequence read length (nt)
E	Last hit Position Query	position on the sequence read corresponding to the farthest HPV16R position to which the sequence read matches
F	Last hit Position HPV	position on the HPV16R corresponding to the farthest position to which the sequence read matches
G	Match nt	number of nucleotides covering the HPV16R hit area
H	Last 2nd-hit Position Query	for sequence with significant second hit, this column contains the position on the sequence read corresponding to the farthest HPV16R position to which the second hit matches
I	Last 2nd-hit Position HPV	for sequence with significant second hit, this column contains the farthest HPV16R position to which the second hit matches
J	Match nt 2ndHit	number of nucleotides covering the HPV16R second hit area
K	Overlap 1st & 2nd hits	for sequences in categories 1.2a and 3.3a, this column contains the number of overlapped nucleotides between the two hits
L	Original Seq Name	original sequence name
M	Reverse complement sequence	reverse complementary sequence of the original sequence read

Program set 3: Multiple-alignment of grouped sequences

This program set performs multiple alignments of sequences in an individual primer group. The program process is illustrated in Figure 2.30. This set consists of two programs, with the following names:

`set3_p8mac.pl`

`set3_p9mac.pl`

The inputs to this program set are FASTA sequence lists of each RA_HP16 primer group of each DNA sample from program set 2. To obtain a good alignment, the sequences should be edited prior to performing an alignment so that they contain as little sequence variations as possible. The program considers only the HPV16 sequence part of each sequence read to generate multiple alignments for each RA_HP16 primer group.

Program set 4: Writing different sequence outputs in FASTA format

Program set 4 re-edits the results of multiple alignments from program set 3. The program process is illustrated in Figure 2.30. It consists of a set of two programs, with the following names:

`set4_p10mac.pl`

`set4_p11mac.pl`

The alignments generated by program set 3 contain only the HPV16 portions of each sequence. This program set re-edits the alignments by restoring the HPV16-portion-only sequences to their original sequences. The edited alignments are then exported in FASTA format as outputs. The alignments can be viewed and edited in any sequence editor software or even a text editor. In this study, Geneious Pro version 4.8.5 was used.

2.2.1.3 Requirements prior to executing the programs

Because the program sets integrate the usages of other stand-alone programs such as BLAST and ClustalW, these stand-alone programs must be installed in the proper locations on the analyzing computer prior to executing the ASP16 analysis program sets. Databases for BLAST analysis must also be set up. The following describes six requirements before the ASP16 analysis programs can be executed. The website addresses where BLAST and ClustalW applications can be obtained are indicated in the Materials and Methods section.

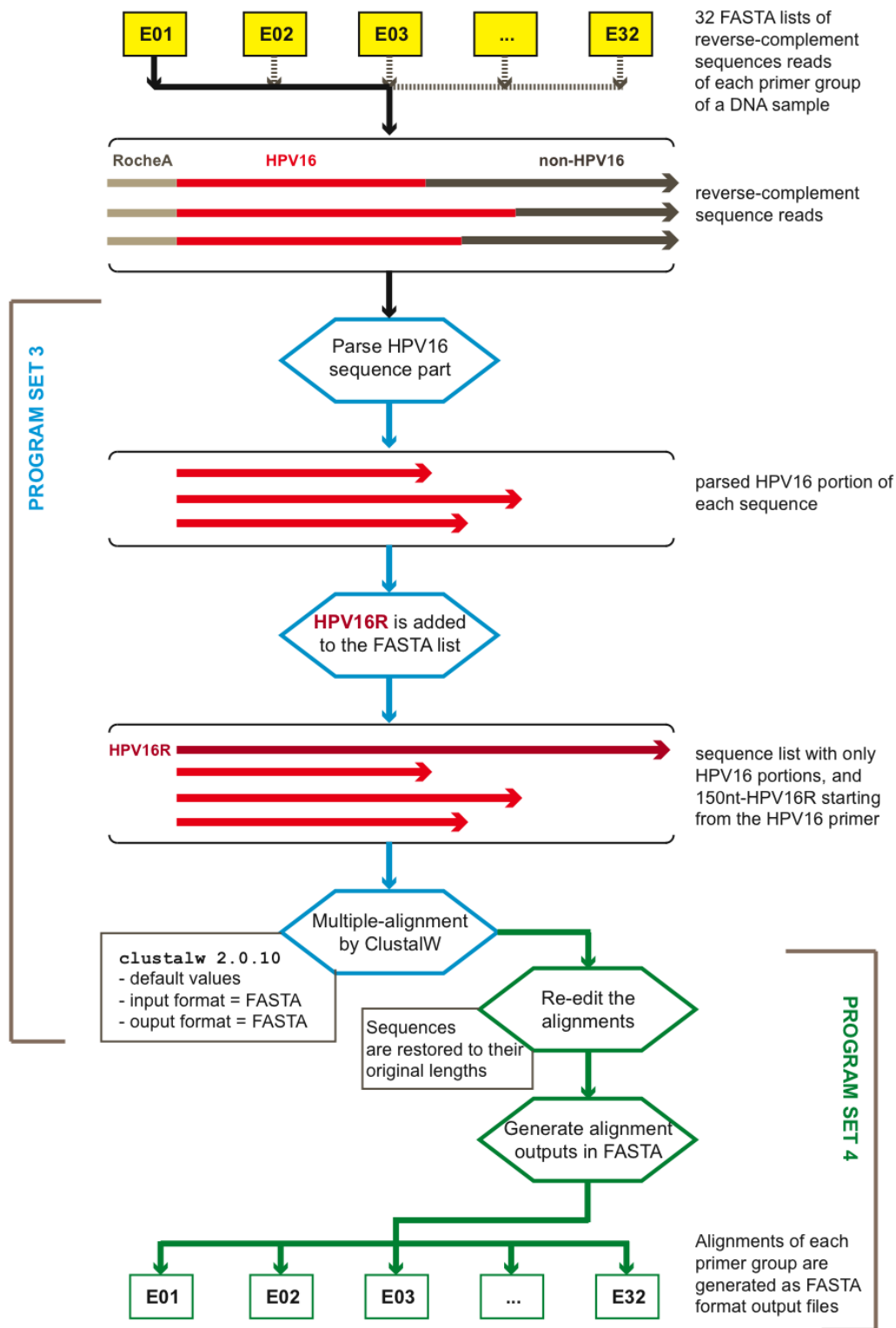


Figure 2.30: Scheme of program sets 3 and 4. In program set 3, the HPV16 portion of the significant sequences are extracted and put in a new list for each RA_HPV16 primer group. A 150-nt HPV16R sequence is then added to the list. The added portion of the HPV16R sequence starts at the same position as the respective RA_HPV16 primer. Multiple alignment is performed to the prepared sequence list, using ClustalW version 2.0.10 with default values. Program set 4 sorts and edits the resulting alignments for simpler sequence visualization. The edited alignments are exported in FASTA format. Each DNA sample group contains thirty-two alignment outputs because 32 primers are used in the multiplex PCR.

(1) Perl program version 5.8.8 or newer must be installed. In Mac OS and Linux OS, Perl is usually already installed. To check if Perl is installed, the command “perl -v” is given in the command-line. It should report the version of the Perl program.

(2) BLAST executables version 2.2.17 (binary for Unix) must be installed under a directory named “blast-2.2.17”, located directly under root.

(3) ClustalW version 2.0.10 for Mac OS must be installed under a directory named “clustalw-2.0.10-macosx”, located directly under root. The text file containing thirty-two 151-nt partial HPV16R nucleotide sequences, named “hvp16R_EEprimer_151nt”, must be saved under the same directory. The file contents are given in Appendix A7.

(4) Databases of the HPV16R reference genome and RA_HP16 primers must be formatted and placed under directory /blast-2.2.17/db/. For the HPV16R database, a text file containing FASTA format of the HPV16R sequence must be created with the name “hvp16”. For the HPV16 primer database, a text file containing FASTA format of the thirty-two RA_HP16 primer sequences (with primer names E01,...,E32) must be created with the name “EEpyro3”. The primer sequences are listed in Materials and Methods section. To format the databases, the stand-alone program from the BLAST executables, `formatdb`, is used with the commands:

“`formatdb -i hvp16 -p F -o T`” for HPV16R database, and

“`formatdb -i EEpyro3 -p F -o T`” for primer database.

After the formatting is completed, five new files are created for each database. These files must be moved under directory /blast-2.2.17/db/.

(5) A text file input must contain total sequence reads from the Roche/454 GS-FLX pyrosequencing, in FASTA format. No fixed path is required for the input file.

(6) The path (or file location) where ASP16 analysis programs are located, and the path where the ASP16 analysis results should locate, must be specified. In each of the four ASP16 program sets (`PROG_SET_1_mac.pl`; `PROG_SET_2_mac.pl`; `PROG_SET_3_mac.pl`; `PROG_SET_4_mac.pl`), the default paths can be changed to user-specified paths by editing the Perl scripts of these programs in a text editor application.

2.2.1.4 Executing the programs

All programs in this study were written in Perl language, and were executed under Perl program version 5.8.8. They had not been tested with other versions of Perl program. The programs were designed for command-line interface, not graphical user interface (GUI).

The GUI is not necessary because the programs were written in such a way that only a few simple short commands are required. This means a user is required to “enter” one or more commands to the computer manually in order to run these ASP16 analysis programs, instead of clicking buttons. In Mac operating system, the command-line interface can be accessed through the Terminal application. To run the complete four program sets, two commands are used.

For program set 1: “perl PROGRAM Input_A Input_C”

For program sets 2 - 4: “perl PROGRAM Input_B Input_C”

where PROGRAM contains complete program name and path
 Input_A is complete path of input FASTA file
 Input_B is barcode number (01, 02, ..., 24)
 Input_C is the experimental number (1, 2, 3, 4, ...)

(1) PROG_SET_1_mac.pl This program sorts sequences into sample groups according to the barcodes. It is required to run only once.

(2) PROG_SET_2_mac.pl This program automatically executes all programs in set 2 (programs 2-7) that perform blasts and statistic calculations. It works for an individual barcode, and is required to run once for each barcode.

(3) PROG_SET_3_mac.pl This program automatically executes all programs in set 3 (programs 8-9) that perform sequence alignments. It should be run only after program set 2 has been completed. It works for an individual barcode, and is required to run once for each barcode.

(4) PROG_SET_4_mac.pl This program automatically executes all programs in set 4 (programs 10-11) that generate full FASTA sequence end-results. Input_B can be any number. It is required to run only once after all barcodes have been analyzed by program set 2 and set 3.

2.2.1.5 Output files used for ASP16 analysis

There are several output files generated by the ASP16 analysis programs. Most output files exist for the benefits of cross-communications between programs. The following output files are informative for the users.

Output 1: Sequences and HPV16R-hit data

These files are generated by program set 2. There is a single file produced per DNA sample. An example of these files is shown in Figure 2.29. The information from

these outputs can be used as reference data sources for the sequence reads. It is also possible to predict a possible HPV16 breakpoint for each DNA sample by sorting/editing the data in Excel.

Output 2: Basic statistics for sequences with or without 28-nt cutoff

Two files are produced per DNA sample, one containing the data without the 28-nt cutoff filtering, and the other with the filtering. These files are generated by program set 2. They contain numbers of sequence reads of each RA_HP16 primer group of each DNA sample. The information was used to perform further statistical analyses of the experimental results, such as calculation of histograms of sequence distribution among different primer group. See Table 2.16 and Table 2.17 for the results of ASP16-3 and ASP16-4.

Output 3: Multiple alignments of sequences in each RA_HP16 primer group

Maximally thirty-two files are generated per DNA sample, corresponding to the 32 RA_HP16 primer groups. The number of the output files depends on the number of the primer combinations used. These files are generated by program set 4. They contain the alignments of sequences in an individual RA_HP16 primer group of a DNA sample in FASTA format. An example of these sequence alignments is shown in Figure 2.39. The alignments can be used to find HPV16R point mutations, to search for HPV16 integration junction sites, and to create an assembled HPV16 sequence for each DNA sample. The assembled sequences can be used to construct a phylogenetic tree. See section 2.2.2 for application and illustration.

2.2.2 Analysis of HPV16 sequences in ASP16-3 and ASP16-4

The first two ASP16 experiments, ASP16-1 and ASP16-2, had been conducted by Bo Xu (Xu, 2010). In these two experiments, the sequence reads of each DNA sample covered about 50% of the analyzed HPV16 E1-E2 gene region, in average. This means that the chance of finding any HPV16 integration site was only about 50%. Further optimization of the ASP16 was required to obtain better sequence coverage.

In the frame of this PhD work, two ASP16 experiments (ASP16-3 and ASP16-4) were performed and the computer programs described in the previous section were used for data analysis. The main purpose was to find HPV16 integration junctions in the analyzed

DNA samples. Another purpose was to continue the optimization of the ASP16 strategy and investigate the results of the optimizations, in particular whether complete sequence coverage could be achieved. In addition, the sequence data were used to identify the HPV16 variants present in the samples. The results of both ASP16 experiments will be described together.

2.2.2.1 Optimization of ASP16

The Roche/454 GS-FLX is capable of amplifying up to 500 bp DNA fragments during the emulsion PCR step, and delivering average sequence reads with length around 200 nt (<http://www.454.com/>). However, the average sequence read lengths obtained in ASP16-1 and ASP16-2 were only 116 nt and 92 nt, respectively (Xu, 2010). Due to the presence of RA_HP16 primer (39-43 nt), GP16 (at least 18 nt) and 4-nt barcode (4 nt) sequences at the termini of each sequence read, the informative sequence of the amplified genomic DNA is at least 61 nt shorter than the total sequence read length (Figure 2.31). Thus, the average informative sequence lengths of ASP16-1 and ASP16-2 were only 55 nt and 31 nt, respectively. Shorter-than-expected average sequence read length is the major reason for low sequence coverage in ASP16-1 and ASP16-2. The problem of short read length is due to the general limitation of the Roche/454 GS-FLX system for amplicon sequencing, not for genomic DNA sequencing (S. Wolf, personal communication).

In ASP16-2, sixteen combinations of HPV16 forward primers were used, spanning ~3000 bp of the HPV16 E1-E2 area (Xu, 2010). Each combination includes a biotin-labeled HPV16 primer and the corresponding RA_HP16 primer. The distance between the primers was about 200 bp. Due to the average sequence read length of ~100 nt obtained by pyrosequencing, about 50% of the E1-E2 sequence region was not covered by the ASP16 sequences (Figure 2.32).

In the ASP16 experiments (ASP16-3 and ASP16-4) performed in this PhD work, an essential aim was to try to obtain 100% sequence coverage in the HPV16 E1-E2 region. Toward this goal, the distance between HPV16 primers was reduced to about 100 bp by adding another sixteen combinations of HPV16 primers. These primer combinations were used since experiment ASP16-3. Two of the additional primer combinations, E17 and E21, replaced two previous combinations of ASP16-2, which were excluded due to their

low specificity. Altogether, 30 combinations of HPV16 primers were used in ASP16-3. The sequences and names of all HPV16 primers for ASP16-3 and ASP16-4 are given in Materials and Methods. The locations of the RA_HPV16 primers of these 30 pairs are shown in Figure 2.33.

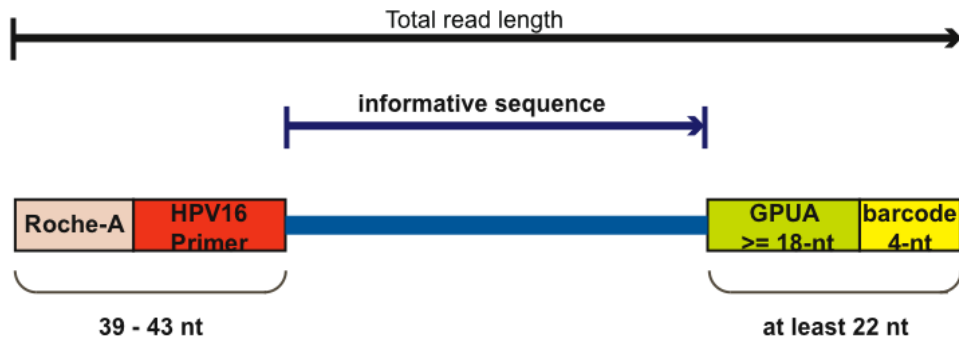


Figure 2.31: Components of an ASP16 sequence read. The total read length of each sequence read contains at least 61 nucleotides derived from the flanking primer sequences. The informative sequence corresponds to the remaining nucleotides of amplified sample DNA.

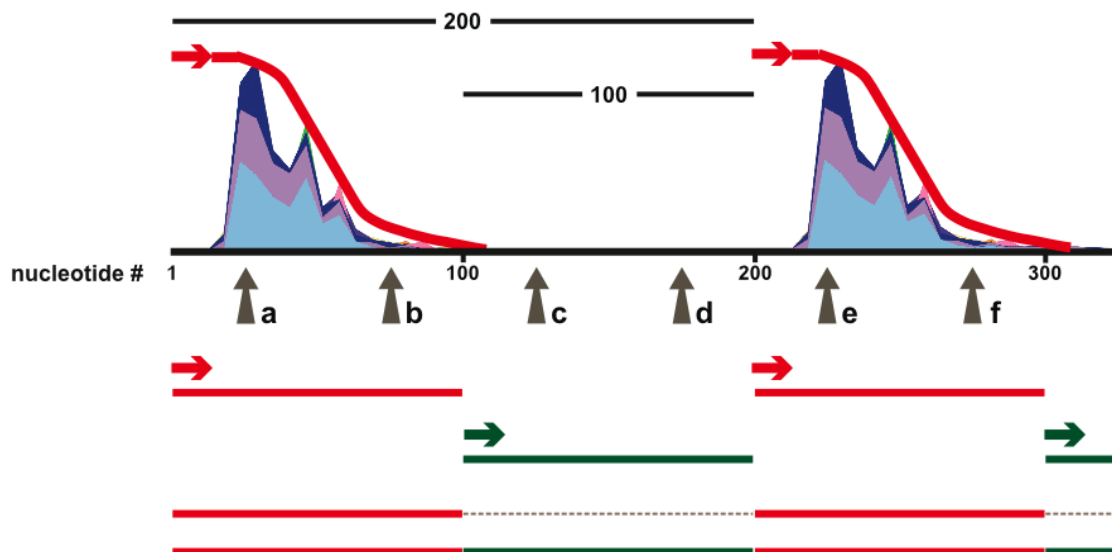


Figure 2.32: Model of primer locations and sequence read length in ASP16 experiments. The graphs represent histograms of sequence read length distribution obtained in ASP16-2. The graphs indicate that most sequence reads extend for 25-50 nt from the primer, and only a few sequences reach 100 nt. The HPV16 primers used in previous experiments (ASP16-1 and ASP16-2) are indicated by red arrows, additional HPV16 primers by green arrows. Red and green lines symbolize sequence reads of 100 nt. The gray vertical arrows (a-f) indicate the positions of possible HPV16 integration breakpoints. Since the sequence read length was ~100 nt but the primers (red arrows) were ~200 bp apart, gaps of about 100 nt with no sequence information were obtained. Addition of new HPV16 primers (green arrows) shortened the average distance between primers to 100 nt, thereby improving the probability for complete sequence coverage.

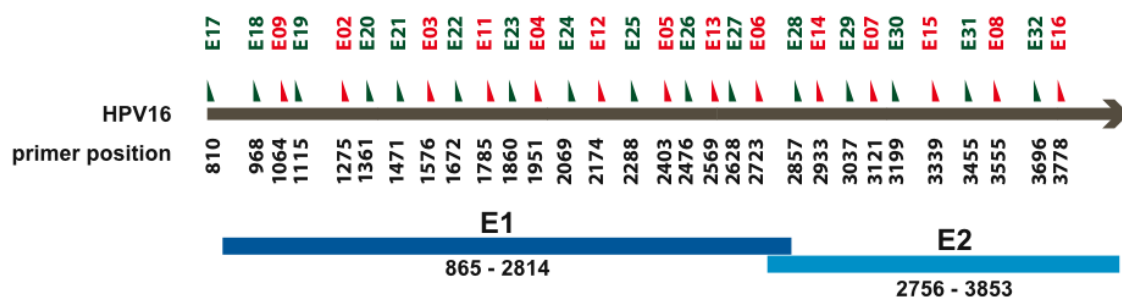


Figure 2.33: Locations of RA_HP16 primers in ASP16-3 and ASP16-4. All thirty RA_HP16 primers used in ASP16-3 and ASP16-4 are shown. All primers are in forward direction and cover almost the complete region of the HPV16 E1 and E2 genes. The red and green triangles indicate positions of previous primers (used in ASP16-2) and additional primers (added since ASP16-3), respectively. The primers were named with numbers according to their entries in the experiment. The positions of ORFs E1 and E2 are indicated at the bottom.

2.2.2.2 DNA samples in ASP16-3 and ASP16-4

Altogether 25 DNA samples were analyzed in ASP16-3 and ASP16-4 (Table 2.13). Four samples were HPV16-positive cervical carcinoma cell lines and 21 samples were cervical scrapes. The four cervical cell lines served as controls because their integration junctions are known. The cervical DNA samples were selected for the experiments based on their high probabilities of having integrated HPV16 (integration %). The integration percentage values were obtained through E2/E6 RT-qPCR (described in Introduction). From the clinical DNA samples, nine were analyzed in both experiments.

2.2.2.3 HPV16 amplicon preparation for ASP16-3 and ASP16-4 sequencing

To prepare HPV16 amplicons for pyrosequencing, the genomic DNA of each sample was subjected to several amplification steps, as described in Introduction and Materials and Methods. In ASP16-3 and ASP16-4, thirty biotin-labeled HPV16 forward primers were used during the linear amplification step, distributed into four primer mixes (Table 2.14). The 7 or 8 HPV16 primers of each primer mix locate at about 400 bp apart on the HPV16 genome. The HPV16 amplicons of each DNA sample were generated by HPV16 semi-nested multiplex PCR using the corresponding bipartite RA_HP16 nested forward primers and a tripartite reverse primer (Roche-B_barcode_GPUA) (Table 2.14, see also Figure 2.33).

Table 2.13: DNA samples in ASP16-3 and ASP16-4.

Sample Name	ASP16-2 ID*	ASP16-3 ID	ASP16-4 ID	Integration % ^(a)	Cytology/histology ^(b)	Origin ^(c)
MRI-H186	2B01	3B01	4B01	49	cancer	cell line
MRI-H196	2B02	3B02	4B02	32	cancer	cell line
SiHa	2B03	3B03	4B03		cancer	cell line
CaSki			4B04		cancer	cell line
1		3B04		76	CIN 2/3	Besancon
2		3B05		74	CIN 2/3	Besancon
5		3B06		95	CIN 2/3	Besancon
1503		3B07	4B12	74	CIN 2/3	Besancon
1511		3B08	4B13	87	CIN 2/3	Besancon
1801		3B09	4B14	72	CIN 2/3	Besancon
2219		3B10	4B15	72	CIN 2/3	Besancon
2227		3B11	4B16	99	CIN 2/3	Besancon
2229		3B12		76	CIN 2/3	Besancon
2237		3B13	4B17	98	CIN 2/3	Besancon
3009		3B14	4B18	81	CIN 2/3	Besancon
3035		3B15		74	CIN 2/3	Besancon
4242		3B16	4B20	79	CIN 2/3	Besancon
66019		3B18	4B06	100	HSIL/CIN3	Reims
07 C 381	2B06		4B05	100	cancer	Reims
61979			4B07	100	HSIL	Reims
75022			4B08	100	LSIL	Reims
75857			4B09	100	HSIL	Reims
4			4B11	82	CIN 2/3	Besancon
4238a			4B19	64	CIN 2/3	Besancon
07 C 368	2B05		4B21	84	cancer	Reims

* ASP16-2 was performed by Bo Xu (Xu, 2010).

(a) The integration percentage was determined by E2/E6 real-time quantitative PCR as described by (Briolat et al, 2007). See Introduction for further explanation.

(b) CIN: cervical intraepithelial neoplasia. HSIL: high-grade squamous intraepithelial lesion.

(c) The clinical DNA samples were obtained from Reims or Besancon.

Table 2.14: Primer combinations for HPV16 amplification in ASP16-3 and ASP16-4.

	HPV16 linear amplification	Semi-nested HPV16 multiplex PCR	
	(FORWARD) Biotin-labeled HPV16 primer ^(a)	(nested FORWARD) RA_HPV16 primer ^(b)	(REVERSE) Roche-B_barcode_GPUA primer ^(c)
Primer mix 1	5B-790F (BE-17)* 5B-1261F (BE-02) 5B-1562F (BE-03) 5B-1938F (BE-04) 5B-2389F (BE-05) 5B-2705F (BE-06) 5B-3101F (BE-07) 5B-3542F (BE-08)	E17 (810)* E02 (1275) E03 (1576) E04 (1951) E05 (2403) E06 (2723) E07 (3121) E08 (3555)	RB-B01, RB-B02, ... or RB-B24
Primer mix 2	5B-1046F (BE-09) 5B-1457F (BE-21) 5B-1760F (BE-11) 5B-2151F (BE-12) 5B-2539F (BE-13) 5B-2912F (BE-14) 5B-3318F (BE-15) 5B-3762F (BE-16)	E09 (1064) E21 (1471) E11 (1785) E12 (2174) E13 (2569) E14 (2933) E15 (3339) E16 (3778)	
Primer mix 3	5B-952F (BE-18) 5B-1336F (BE-20) 5B-1849F (BE-23) 5B-2277F (BE-25) 5B-2613F (BE-27) 5B-3021F (BE-29) 5B-3444F (BE-31)	E18 (968) E20 (1361) E23 (1860) E25 (2288) E27 (2628) E29 (3037) E31 (3455)	
Primer mix 4	5B-1105F (BE-19) 5B-1653F (BE-22) 5B-2053F (BE-24) 5B-2461F (BE-26) 5B-2842F (BE-28) 5B-3189F (BE-30) 5B-3680F (BE-32)	E19 (1115) E22 (1672) E24 (2069) E26 (2476) E28 (2857) E30 (3199) E32 (3696)	

* Primers BE-17 and E17 were not used in ASP16-4

(a) BE = biotin-labeled HPV16 primers.

(b) E = RA_HPV16 primers. The numbers in brackets are the positions (5'end) of the primers on HPV16R genome.

(c) RB = Roche-B_barcode_GPUA primers.

In ASP16-3, the HPV16 amplicons were analyzed by agarose gel electrophoresis and Southern hybridization, to determine whether the amplicons were HPV16-specific. The agarose gels and hybridization results are shown in Figure 2.34. The HPV16 PCR products were expected to contain fragments of variable sizes, producing smears on the gel. This expectation was fulfilled in the majority of cases. Some amplicons showed prominent bands (such as panel B, lane 9-3 or panel D, lane 13-2). These bands were not HPV16-specific, as no signal was observed after hybridization. The majority of the amplicons contained HPV16-specific fragments distributed between 100-600 bp. Many amplicons contained even larger fragments with sizes up to 1.5-2.0 kb. Aliquots of 100 ng from each amplicon were pooled together for Roche/454 GS-FLX sequencing. In ASP16-4, each HPV16 amplicon was size-selected to 200-300 bp using E-Gel SizeSelect system (Invitrogen), before being pooled together for sequencing, as described in Materials and Methods.

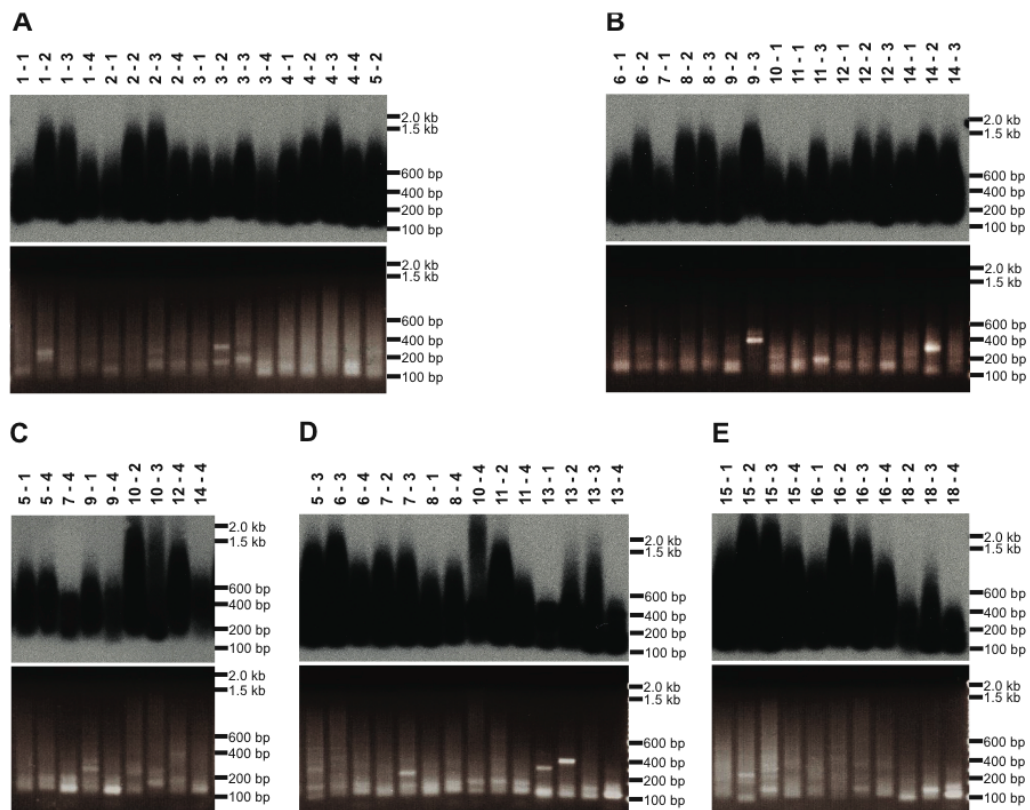


Figure 2.34: HPV16 multiplex PCR products in ASP16-3. Each panel (A-E) shows an agarose gel photo in the lower part and the Southern blot in the upper part. 50-ng aliquots of HPV16 multiplex PCR amplicons were loaded on the gel. The labeling of each lane indicates the barcode number (1, 2, 3,..., 24), followed by the multiplex reaction mix number (1, 2, 3, 4). The gels were blotted and hybridized with complete HPV16 genome as ^{32}P -labelled probe.

2.2.2.4 Amplicon sequencing by Roche/454 GS-FLX pyrosequencing

The pooled HPV16 amplicons were sequenced, using the Roche/454 GS-FLX standard system, at the DKFZ Genomics and Proteomics Core Facilities. GS emPCR kit III was used for emulsion PCR and GS LR70 as the sequencing kit, with one large region on a 70x75 picotiter plate used. Roche-B was used as the sequencing primer. The manufacturer's estimated number of sequence reads per sequencing run for the described sequencing format is 210000. The results were delivered as FASTA files, containing complete sequence reads of each experiment. The basic features of both experiments are shown in Table 2.15, in comparison with the results of the previous experiment ASP16-2 (Xu, 2010).

Table 2.15: Features of ASP16-2, ASP16-3 and ASP16-4

	ASP16-2*	ASP16-3	ASP16-4
Number of DNA samples	19	17	20
Number of HPV16 primer combinations	16	30	29
Number of total sequence reads	220001	152147	129150
Average sequence read length, excluding 19-nt Roche-B sequence (nt)	92	105	108
Average number of sequence reads per primer group per DNA sample	724	298	222

* Experiment ASP16-2 was conducted by Bo Xu (Xu, 2010).

2.2.2.5 Statistics of ASP16-3 and ASP16-4

The FASTA files containing the total sequence reads of each experiment, were initially analyzed with the four ASP16 analysis program sets, described in section 2.2.1.4. The numbers of sequence reads per primer group per sample group, before and after filtering with 28-nt cutoff, are shown in Table 2.16 for ASP16-3 and in Table 2.17 for ASP16-4.

To evaluate whether the HPV16 primers produced HPV16-specific sequences, the percentage of significant HPV16 sequence reads for each sample and each primer group was calculated with the following formula:

$$\frac{(\text{number_of_sequence_reads_after_28nt_cutoff})}{(\text{number_of_sequence_reads_before_28nt_cutoff})} \times 100\%$$

For each primer group, the distribution of the percentages of significant HPV16 sequences from all DNA sample are shown in boxplot diagrams (Figure 2.35). Twenty-eight of the thirty primers showed high specificity of more than 50%, which means that more than 50% of the sequence reads in the primer groups were HPV16-specific. Primer E17 (810) showed a reduced efficiency. Because it is located in the 5' part of the E1 gene where integration is unlikely to occur, it was not used in ASP16-4 and no alternative primer was added to replace primer E17 (810). Primer E21 (1471) also showed reduced efficiency, but was included also in ASP16-4.

Table 2.16: Number of sequence reads per primer group of ASP16-3 before and after 28-nt cutoff.

NO	CF	810	968	1064	1115	1275	1361	1471	1576	1672	1785	1860	1951	2069	2174	2288	2403	2476	2569	2628	2723	2857	2933	3037	3121	3199	3339	3455	3555	3696	3778	Sum	Total
		E17	E18	E09	E19	E02	E20	E21	E03	E22	E11	E23	E04	E24	E12	E25	E05	E26	E13	E27	E06	E28	E14	E29	E07	E30	E15	E31	E08	E32	E16	sigEE	HPV
3B01	7	294	646	295	521	287	127	554	240	98	308	350	255	71	38	15	42	71	79	214	123	9	92	263	184	279	195	438	328	122	6368	6587	
3B02	13	271	629	346	265	147	153	721	222	100	314	295	336	37	43	9	42	27	43	66	109	31	111	376	279	360	247	652	517	420	6936	7116	
3B03	10	137	500	220	167	90	195	512	283	104	423	132	553	135	254	19	56	23	194	109	149	37	386	28	32	51	179	498	529	200	6205	6361	
3B04	81	294	1045	378	444	165	389	91	486	201	391	408	554	284	237	30	81	86	140	342	233	99	378	538	495	713	266	784	694	681	11827	12144	
3B05	12	189	1043	332	401	58	337	913	425	149	451	283	605	251	188	18	56	45	117	161	140	74	298	254	386	659	234	887	600	720	10326	10612	
3B06	8	214	1250	418	885	69	157	1156	362	141	444	343	518	117	70	47	67	30	18	232	13	60	188	684	551	333	333	692	134	110	9642	9870	
3B07	63	504	975	420	592	148	259	1101	387	173	611	345	913	130	206	71	121	67	142	234	225	33	218	336	380	834	374	1009	498	545	11914	12166	
3B08	19	538	1395	392	325	160	195	923	328	196	820	177	470	172	289	72	173	86	75	255	9	18	433	406	602	434	408	432	36	94	9934	10173	
3B09	17	61	829	197	71	16	163	780	373	111	396	70	740	117	234	16	75	47	141	92	194	27	254	184	170	650	326	772	136	455	7624	7856	
3B10	55	423	980	176	453	78	23	1559	248	80	188	131	297	37	100	27	128	13	8	141	24	1	2	115	205	390	634	662	106	2	7286	7572	
3B11	19	374	832	271	536	206	161	1136	277	125	811	306	649	100	148	23	73	26	138	266	218	11	370	471	400	477	457	1157	482	398	10918	11194	
3B12	2	376	918	323	700	361	340	1135	450	207	596	223	554	337	101	28	146	38	176	230	144	24	218	330	492	676	561	1092	306	588	11672	12013	
3B13	70	367	1210	342	90	50	44	1188	266	59	768	205	804	18	385	4	81	24	30	119	6	21	22	106	598	208	647	202	21	57	8012	8221	
3B14	9	61	890	426	223	1	104	720	276	72	396	156	684	112	14	7	50	16	8	114	20	82	230	272	596	324	309	473	93	148	6886	7073	
3B15	2	438	591	400	359	11	64	832	274	91	872	229	543	124	21	19	86	58	31	214	70	208	170	288	527	326	541	775	244	358	8766	9025	
3B16	1	216	460	284	397	29	82	562	135	242	373	213	514	145	3	5	55	16	13	106	16	49	46	135	485	266	313	555	53	157	5926	6093	
3B18	1	631	1202	167	3	160	193	2	280	32	376	0	1285	26	303	0	28	9	17	1	1	9	17	1	143	86	28	3	28	60	5092	5237	

B

CF	E=8	810	968	1064	1115	1275	1361	1471	1576	1672	1785	1860	1951	2069	2174	2288	2403	2476	2569	2628	2723	2857	2933	3037	3121	3199	3339	3455	3555	3696	3778	Sum	Total
		E17	E18	E09	E19	E02	E20	E21	E03	E22	E11	E23	E04	E24	E12	E25	E05	E26	E13	E27	E06	E28	E14	E29	E07	E30	E15	E31	E08	E32	E16	sigEE	HPV
3B01	7	294	338	291	489	287	52	519	240	93	204	327	252	66	30	15	37	64	69	186	117	9	89	172	184	204	131	312	270	97	5445	6587	
3B02	9	270	461	341	23	147	55	606	219	85	218	1283	334	37	33	6	40	25	30	53	107	30	103	259	278	289	155	484	463	328	5771	7116	
3B03	2	125	338	206	160	81	64	408	259	68	283	285	337	119	47	13	26	20	143	91	138	32	357	1	0	0	4	381	389	134	4507	6361	
3B04	17	270	681	364	417	151	102	676	467	159	177	364	545	246	32	13	44	67	100	273	208	80	340	287	470	487	158	508	605	439	8747	12149	
3B05	5	182	638	323	385	55	91	741	412	107	215	248	593	198	35	15	34	29	85	144	127	59	272	133	364	433	122	599	520	406	7570	10612	
3B06	6	214	885	397	852	66	49	993	358	118	264	312	506	109	56	37	49	25	12	200	13	54	164	500	548	256	224	558	118	66	8009	9870	
3B07	10	483	613	348	516	134	112	778	337	126	374	255	885	106	130	50	69	51	114	190	203	27	206	188	314	546	197	642	400	363	8767	12166	
3B08	7	531	939	362	299	159	87	691	319	158	543	137	455	159	233	46	132	80	62	191	9	18	379	277	574	384	228	349	19	42	7869	10173	
3B09	5	55	458	151	57	15	56	576	358	75	256	53	712	93	64	11	38	40	106	72	177	18	234	113	166	275	114	506	112	295	5261	7856	
3B10	0	421	844	139	445	76	18	803	174	74	138	80	230	15	86	19	40	13	8	31	7	1	2	2	45	383	325	282	12	0	4713	7572	
3B11	8	365	562	243	475	198	58	925	274	89	495	273	641	85	133	16	45	16	110	214	206	11	336	365	377	353	245	849	373	278	8618	11194	
3B12	0	375	697	314	668	359	222	875	446	195	398	194	550	320	88	22	83	32	144	203	141	24	206	232	492	500	358	768	259	429	9594	12013	
3B13	0	332	718	324	81	19	11	581	255	45	480	186	801	9	3	2	36	19	22	16	2	10	12	57	550	54	347	105	15	25	5117	8221	
3B14	6	61	603	393	209	1	24	607	276	52	258	140	671	96	8	3	27	11	5	98	20	72	219	212	595	294	203	393	87	114	5758	7073	
3B15	2	435	454	387	351	11	23	646	272	81	573	221	542	122	18	19	74	58	27	195	69	186	157	225	526	309	383	655	237	209	7467	9025	
3B16	0	213	320	274	390	29	29	478	135	207	237	192	505	139	2	5	39	16	12	95	16	47	42	124	482	221	229	467	53	102	5100	6093	
3B18	1	615	525	162	3	139	13	1	275	14	273	0	1198	12	84	0	14	0	8	1	0	2	6	1	1	0	3	3	0	19	3373	5233	

Panel A shows the numbers before the 28-nt cutoff filtering was applied.

Panel B shows the numbers after the 28-nt cutoff filtering was applied.

The column Sum-sigEE-sigHPV contains sums of sequence reads per sample with significant primer hit.

The column Total-per-sample shows the sum of total sequence reads per sample before any selection process.

Gray areas indicate values excluded from analysis.

Table 2.17: Number of sequence reads per primer group of ASP16-4 before and after 28-nt cutoff.

NO CF																																		Sum sigEE sigHPV	Total per sample
	810	968	1064	1115	1275	1361	1471	1576	1672	1785	1860	1951	2069	2174	2288	2403	2476	2569	2628	2723	2857	2933	3037	3121	3199	3339	3455	3555	3696	3778					
	E17	E18	E09	E19	E02	E20	E21	E03	E22	E11	E23	E04	E24	E12	E25	E05	E26	E13	E27	E06	E28	E14	E29	E07	E30	E15	E31	E08	E32	E16					
4B01	0	208	404	495	89	77	39	645	259	36	441	329	271	119	45	39	78	19	59	540	67	107	95	159	161	163	360	382	408	132	6226	6620			
4B02	0	259	732	399	13	58	90	936	303	113	542	259	715	118	96	27	44	26	83	206	115	97	231	226	463	652	855	403	535	227	8823	9456			
4B03	0	410	657	631	189	219	34	400	303	116	609	383	546	241	248	16	64	31	94	239	93	110	258	1	44	6	6	312	507	93	6860	7411			
4B04	0	255	602	1972	134	182	46	499	555	52	288	188	396	97	34	34	252	22	23	172	156	50	177	313	930	297	382	324	856	278	9566	10491			
4B05	0	387	723	288	168	133	40	641	134	125	640	131	253	149	177	35	23	44	108	53	1	2	3	0	3	15	5	2	1	2	4286	4581			
4B06	0	165	524	76	4	139	78	230	4	28	92	6	11	44	344	2	13	7	9	9	0	4	6	3	7	91	19	13	1	19	1948	2191			
4B07	0	279	272	90	3	59	30	357	1	44	65	2	19	80	320	0	6	3	4	7	1	4	0	8	4	49	50	3	0	8	1768	1982			
4B08	0	81	472	16	26	77	83	151	20	21	130	2	22	12	534	3	0	17	8	6	2	19	41	9	25	155	140	42	4	27	2145	2392			
4B09	0	141	452	1479	51	18	21	455	598	57	369	213	368	185	52	33	123	62	48	219	112	91	258	279	606	670	516	756	987	247	9466	10049			
4B11	0	195	572	720	164	77	44	677	372	66	428	206	363	97	45	19	51	13	40	150	48	86	157	209	458	564	556	509	612	238	7736	8291			
4B12	2	423	688	535	126	175	107	629	340	122	670	140	321	158	209	39	97	24	122	106	69	36	174	192	432	675	525	357	411	216	8120	8755			
4B13	0	298	817	1216	125	15	61	654	509	94	522	145	414	144	268	83	149	48	40	158	7	20	216	242	902	538	469	336	62	50	8602	9344			
4B14	0	137	392	249	50	35	131	626	199	85	411	118	232	110	88	28	35	28	64	114	38	19	124	124	254	484	416	287	228	168	5274	5618			
4B15	0	167	540	114	130	27	103	261	225	62	93	4	171	91	207	11	36	44	12	27	11	30	10	72	180	227	256	97	38	126	3372	3863			
4B16	0	48	500	364	47	12	84	337	415	22	469	41	280	143	79	29	79	30	93	86	101	31	207	133	435	231	389	365	526	197	5773	6105			
4B17	0	86	351	114	32	24	66	286	62	58	159	32	66	36	310	7	17	13	17	22	2	22	22	76	119	165	209	117	11	28	2529	2960			
4B18	0	122	643	681	8	0	25	132	277	43	387	15	358	325	62	52	72	88	23	330	11	260	248	417	421	251	394	425	246	260	6576	7123			
4B19	0	72	397	334	24	16	22	261	302	14	481	53	206	31	28	23	68	14	14	203	62	108	276	428	255	364	494	479	375	336	5740	6107			
4B20	0	127	352	840	30	6	46	192	372	21	196	6	271	156	9	18	54	54	13	200	23	193	262	257	838	164	546	288	232	137	5903	6404			
4B21	0	458	584	402	137	65	71	1017	232	96	648	219	207	57	17	8	8	25	38	59	36	26	125	103	310	558	238	385	290	101	6520	7074			

CF #28																																		Sum sigEE sigHPV	Total per sample
	810	968	1064	1115	1275	1361	1471	1576	1672	1785	1860	1951	2069	2174	2288	2403	2476	2569	2628	2723	2857	2933	3037	3121	3199	3339	3455	3555	3696	3778					
	E17	E18	E09	E19	E02	E20	E21	E03	E22	E11	E23	E04	E24	E12	E25	E05	E26	E13	E27	E06	E28	E14	E29	E07	E30	E15	E31	E08	E32	E16					
4B01	0	207	264	495	84	77	14	615	258	33	309	324	264	115	27	39	69	18	55	493	66	102	90	121	161	151	290	311	384	118	5554	6620			
4B02	0	259	422	398	10	57	65	859	299	100	368	244	713	105	74	25	37	23	81	201	109	91	228	171	463	607	718	283	479	202	7691	9456			
4B03	0	410	435	626	185	217	15	372	302	110	489	380	542	229	226	15	61	31	94	238	84	103	234	1	42	0	4	249	460	87	6241	7411			
4B04	0	255	378	1971	130	181	23	451	554	47	214	181	390	85	32	32	223	21	22	170	153	44	171	236	930	276	295	234	829	254	8782	10491			
4B05	0	382	456	287	152	130	23	582	132	122	481	125	246	131	112	27	18	23	95	48	0	0	1	0	0	5	0	0	0	0	3578	4581			
4B06	0	133	191	76	0	109	3	57	1	2	23	5	5	44	30	2	12	0	1	0	0	0	1	0	0	2	0	0	0	0	697	2191			
4B07	0	248	37	90	0	52	5	157	1	9	36	0	17	45	55	0	6	0	0	0	1	0	0	0	0	0	1	0	0	0	760	1982			
4B08	0	14	5	16	3	31	8	63	13	0	14	0	1	0	3	0	0	0	0	0	1	0	0	0	0	0	1	34	3	0	2	212	2392		
4B09	0	139	278	1470	46	16	7	384	597	44	239	197	360	151	42	27	108	55	43	217	107	81	251	206	603	583	364	531	954	197	8297	10049			
4B11	0	194	331	717	157	76	24	591	372	61	291	189	361	87	30	18	48	9	33	148	46	76	145	156	456	524	427	332	571	206	6676	8291			
4B12	2	413	396	528	112	167	65	492	335	110	443	138	315	149	159	34	90	22	106	96	62	32	166	130	431	600	361	251	371	189	6765	8755			
4B13	0	288	445	1206	92	14	43	579	504	78	333	129	405	111	241	68	133	46	36	152	7	13	208	187	900	519	351	261	60	36	7445	9344			
4B14	0	133	185	245	43	31	83	550	191	74	275	112	227	95	64	25	29	24	57	102	35	19	121	87	247	348	281	185	223	130	4221	5618			
4B15	0	142	136	74	111	12	21	160	159	23	51	4	144	11	45	4	16	12	6	23	7	3	2	4	106	96	118	21	2	0	1513	3863			
4B16	0	48	289	362	45	12	56	288	410	17	265	29	272	126	52	26	41	29	83	76	96	25	194	98	426	164	275	272	481	166	4723	6105			
4B17	0	70	85	113	9	6	3	77	36	10	66	20	59	24	5	0	6	3	2	15	2	3	5	6	75	72	161	4	3	1	941	2960			
4B18	0	121	410	678	8	0	14	120	277	30	243	14	354	276	46	50	62	82	17	313	11	227	243	346	418	239	316	334	244	208	5701	7123			
4B19	0	72	235	330	21	16	3	231	301	12	238	47	205	27	23	20	57	14	11	197	59	97	264	275	253	328	362	325	351	277	4651	6107			
4B20	0	127	198	839	26	5	24	181	368	18	122	6	265	142	4	17	47	52	12	189	22	182	246	216	832	159	435	227	226	113	5300	6404			
4B21	0	458	357	385	121	62	47	925	231	85	424	219	206	49	9	7	8	24	33	59	32	25	123	96	310	522	194	328	290	63	5692	7074			

Panel A shows the numbers before the 28-nt cutoff filtering was applied.

Panel B shows the numbers after the 28-nt cutoff filtering was applied.

The column Sum-sigEE-sigHPV contains sums of sequence reads per sample with significant primer hit.

The column Total-per-sample shows the sum of total sequence reads per sample before any selection process.

To determine the distribution of sequence read lengths, histograms of sequence read lengths were plotted for ASP16-2, ASP16-3 and ASP16-4 (Figure 2.36). There were four apparent intervals with peaks at 55, 75, 95 and 150 nt, where most of the sequences fit. For ASP16-2, most sequence reads in the 55-nt peak belonged to primer-dimer products of HPV16 primer E10 (Xu, 2010). This primer was omitted in ASP16-3 and ASP16-4. Despite lower sequence read numbers in ASP16-3 and ASP16-4 compared to ASP16-2 (Table 2.15), the histograms show that higher numbers of long sequence reads (150 nt and longer) have been obtained

indicates that the alterations of linear amplification reaction condition by increasing the elongation time from 30 to 40 seconds have been successful (see Materials and Methods).

To evaluate whether the sequence reads in each size-interval were HPV16-specific, histograms of sequence read lengths before and after 28-nt cutoff filtering were plotted (Figure 2.37). The histograms show that the 75-nt peaks of ASP16-3 and ASP16-4 contain large portions of sequences that are not HPV16-specific. Nevertheless, these sequences are too short to be informative. In contrast, the longer sequence reads (from 95 nt) of both ASP16-3 and ASP16-4 were almost 100% HPV16-specific, and therefore contained informative sequence data.

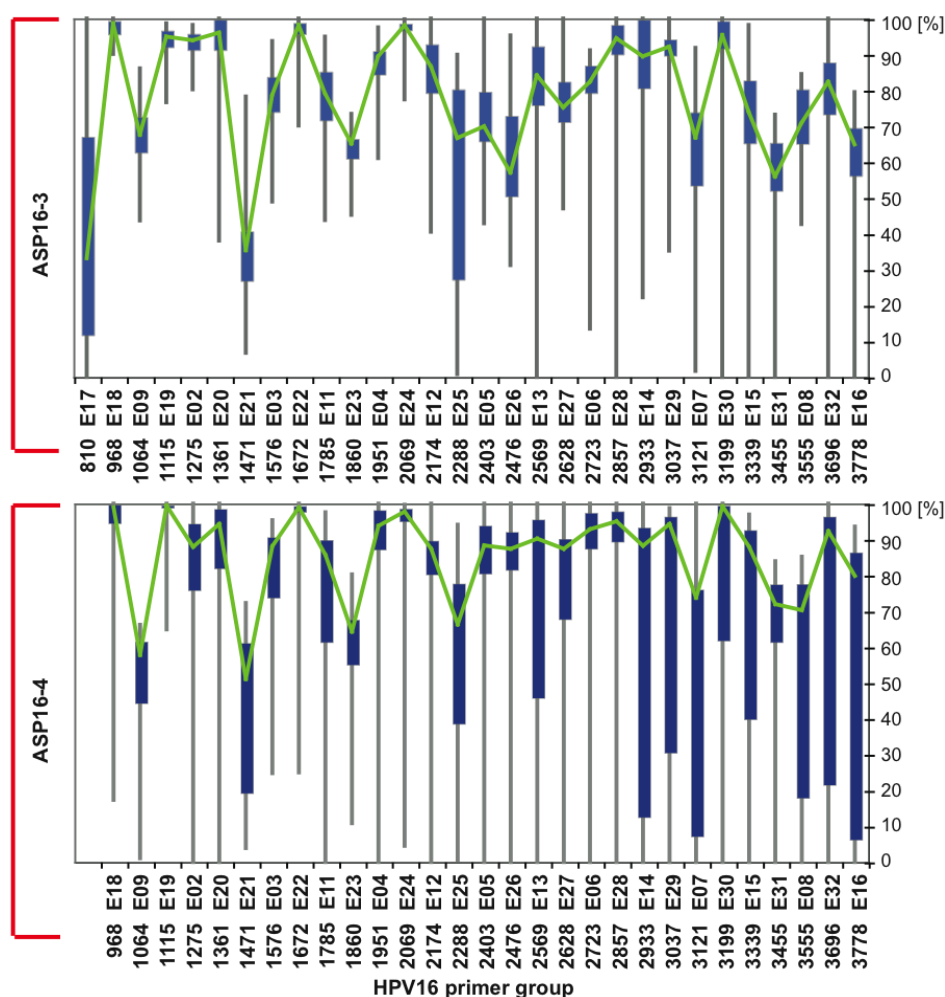


Figure 2.35: HPV16 primer efficiency in ASP16-3 and ASP16-4. Each boxplot shows the distribution (in percentage) of significant HPV16 sequence reads for each HPV16 primer group of all DNA samples. Each blue box contains values from the 25th to 75th percentile. The lines above and below the box represent values above the 75th and below the 25th percentile, respectively. Green lines were plotted from median percentage values of each primer group.

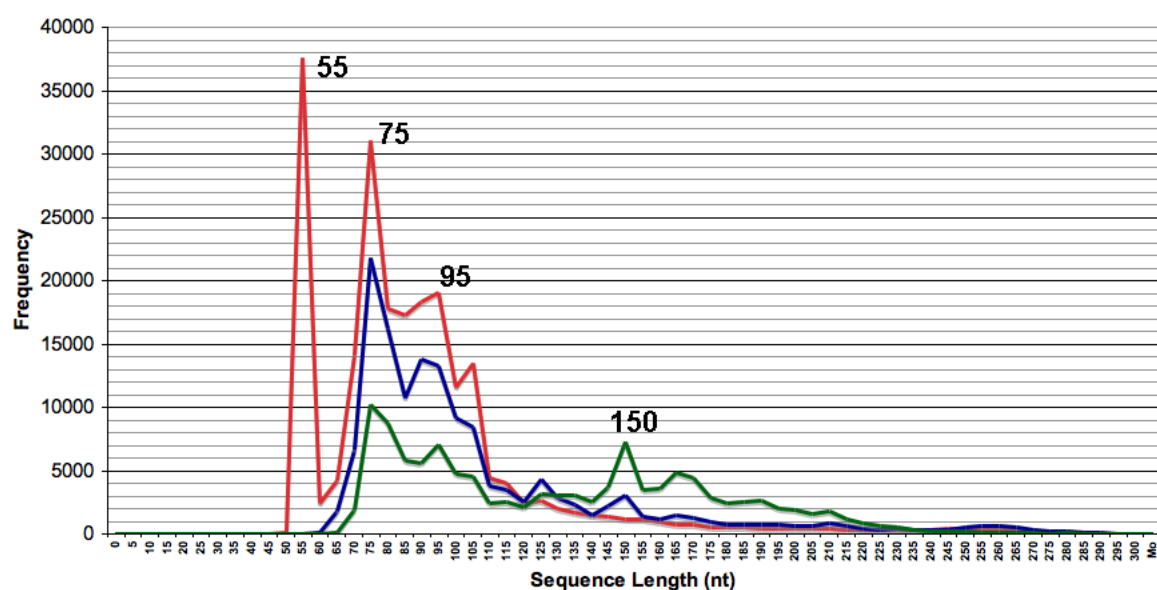


Figure 2.36: Sequence read lengths of ASP16-2, ASP16-3 and ASP16-4. The histograms of sequence read lengths before filtering by 28-nt cutoff are shown. X-axis represents 5-nt intervals of sequence read length. Y-axis indicates the frequency with which sequence read lengths fall into the interval categories. The red line represents ASP16-2, the blue line ASP16-3 and the green line ASP16-4.

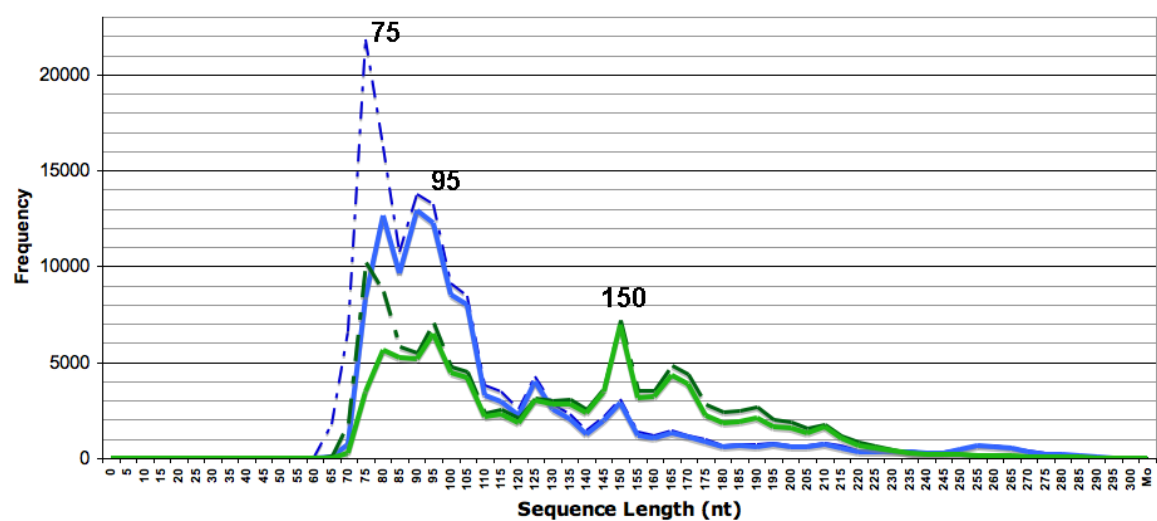


Figure 2.37: Sequence read lengths of ASP16-3 and ASP16-4 before and after filtering with a 28-nt cutoff. The histograms of sequence read lengths before and after filtering by 28-nt cutoff are shown by dashed and solid lines, respectively. X-axis represents 5-nt intervals of sequence read length. Y-axis indicates the frequency with which sequence read lengths fall into the interval categories. The blue lines represent ASP16-3 and the green lines ASP16-4.

2.2.2.6 HPV16 integration junctions

To determine whether the sequences contain possible HPV16 integration junctions, two output data after the ASP16 analysis programs were analyzed.

As first approximate method, the numbers of reads per each primer group of each sample after 28-nt cutoff (shown in Table 2.16 and Table 2.17, both in panel B) were inspected. If the sample contains exclusively integrated HPV16 DNA, the number of reads of the primer groups at and after the HPV16 breakpoint should drop to zero. As an example, the results for the cell line SiHa are shown in Figure 2.38. For the integrated HPV16 genome in SiHa, it is known that the 5' junction is located at pos. 3385 and the 3' junction at pos. 3133. Further more, a small region from pos. 3460-3512 is deleted (Meissner, 1999). In complete agreement with this structure, the read numbers are 1 or zero for primers E07 (3121), E30 (3199), E15 (3339) and E31 (3455) (Figure 2.38).

As the main method, the sequence alignments of each primer group for each DNA sample were analyzed. In this way, the nucleotide sequences at the HPV16 breakpoint location were determined. The alignments were viewed and edited in Geneious Pro version 4.8.5. For illustration, the results for SiHa are shown in Figure 2.39.

Applying these two methods, HPV16 integration junctions were identified in altogether 9 samples in ASP16-3 and ASP16-4: 3 of 4 cell lines, and 6 of 21 clinical samples (Table 2.18). In case of newly identified HPV16 integration junctions, they were examined by junction-specific PCR (Figure 2.40). The positive PCR products were cloned and sequenced. If the cloned sequences confirmed the ASP16 sequences, the junction was declared to be genuine. The next parts describe the determination of HPV16 integrations for all DNA samples listed in Table 2.13, compare the results of the different ASP16 experiments, and examine the ASP16 integration results in relation to those of E2/E6 qPCR.

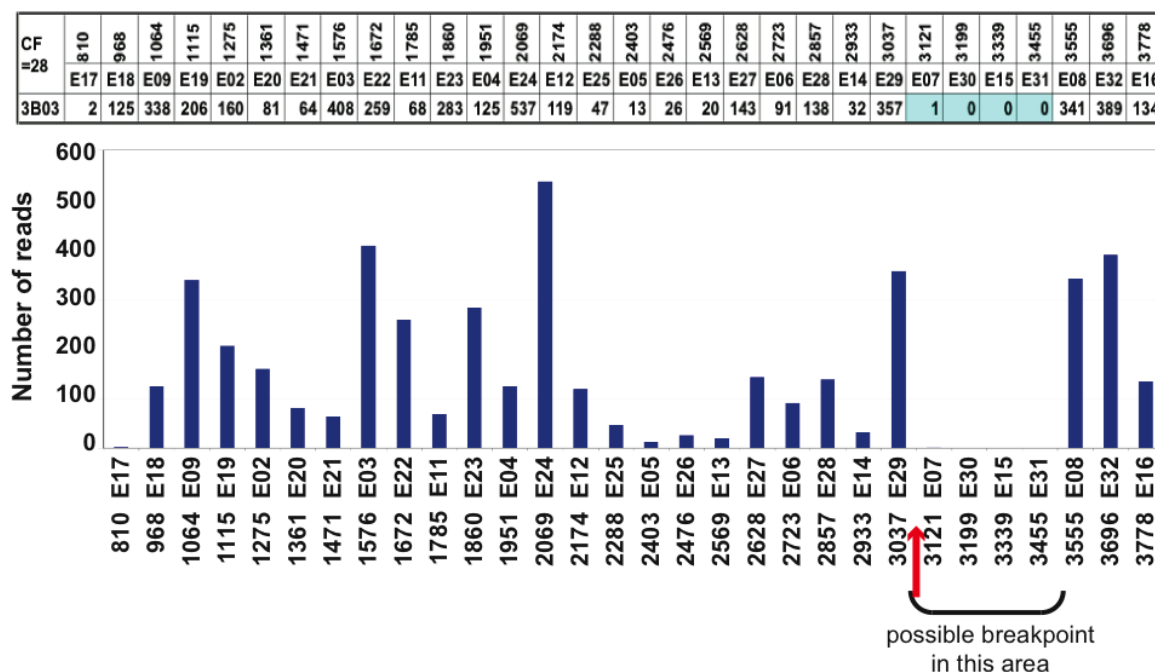


Figure 2.38: HPV16 breakpoint location in SiHa predicted by number of sequence reads. The numbers of reads of each primer group of SiHa in ASP16-3 after 28-nt cutoff, taken from Table 2.16, are shown at the top. A graph was plotted from these values, with X-axis showing the primer groups and Y-axis the number of reads. The possible breakpoint area is shown. The actual HPV16 breakpoint of SiHa is at position 3133 (red arrow).

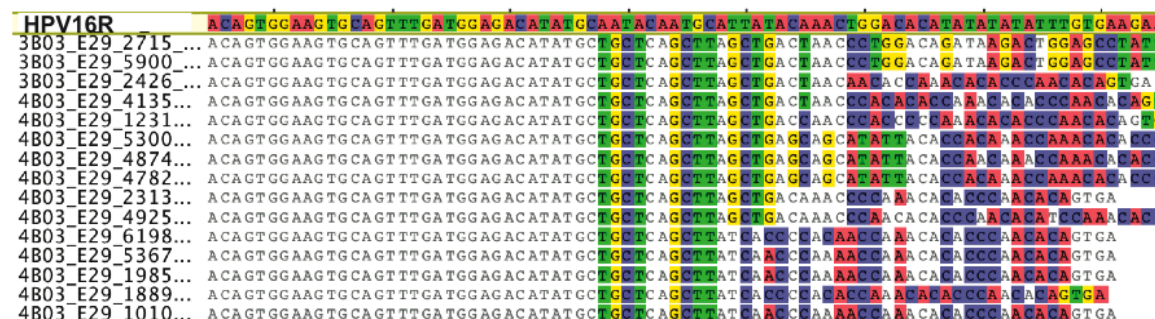


Figure 2.39: Sequence alignment of primer group E29 (3037) of SiHa with the HPV16 integration junction. Bases are color-coded: A in red, C in blue, G in yellow and T in green. The names of the sequences are indicated on the left. HPV16R (top) was used as the reference sequence and every base was highlighted. The sequences in primer group E29 of SiHa in this alignment were taken from both ASP16-3 (3B03) and ASP16-4 (4B03) as the names indicate. The bases in SiHa sequences are highlighted only if they are different from the reference. The highlighted region indicates cellular DNA downstream of the HPV16 integration breakpoint at pos. 3133.

Table 2.18: Identification of HPV16 integration junctions in ASP16 experiments.

Sample name ^(a)	Previous identification: Reference (Method)	Identification in ASP16-2 ^(b)	Identification in ASP16-3	Identification in ASP16-4	Integration percentage by E2/E6 qPCR ^(c)
Cell lines					
MRI-H186	Xu, 2010 (RS-PCR)	yes	yes	yes	49%
MRI-H196	Xu, 2010 (RS-PCR)	(yes) ^(d)	yes	yes	32%
SiHa	Baker, 1987 (Cloning, sequencing)	no	yes	yes	100% ^(f)
CaSki	Smits, 1991 (cDNA sequencing)	no		no	0% ^(f)
Clinical samples for which HPV16 integration breakpoints were identified by ASP16					
CA-07C381		yes		yes	100%
CA-07C368		yes ^(e)		yes	84%
HSIL-66019			yes	yes	100%
HSIL-61979				yes	100%
HSIL-75857				yes	100%
CIN2/3-1801			yes	yes	72%
Clinical samples for which HPV16 integration breakpoints were not identified by ASP16*					
CIN2/3-4242			false-positive	no	79%
CIN2/3-1503			false-positive	no	74%
LSIL-75022				no	100%
CIN2/3-2227			no	no	99%
CIN2/3-2237			no	no	98%
CIN2/3-0005			no		95%
CIN2/3-1511			no	no	87%
CIN2/3-0004				no	82%
CIN2/3-3009			no	no	81%
CIN2/3-0001			no		76%
CIN2/3-2229			no		76%
CIN2/3-0002			no		74%
CIN2/3-3035			no		74%
CIN2/3-2219			no	no	72%
CIN2/3-4238a				no	64%

(a) The sample names for cervical lesions start with the cytological/histological status. CA: cancer. CIN: cervical intraepithelial neoplasia. HSIL: high-grade squamous intraepithelial lesion. LSIL: low-grade squamous intraepithelial lesion.

(b) ASP16-2 was performed by Bo Xu (Xu, 2010).

(c) See Introduction for details of E2/E6 qPCR.

(d) One ASP16 read contained only 6 bp of cellular sequences, too short to be identified by Blast. The sequence was identified as integration junction by comparison with the RS-PCR sequence.

(e) Artifact integration junction was identified from the ASP16 reads. The genuine junction was identified via cloning and sequencing, see text for details.

(f) Estimated integration percentages.

* Ordering of samples by decreasing E2/E6 qPCR values for integration, except for samples CIN2/3-4242 and CIN2/3-1503.

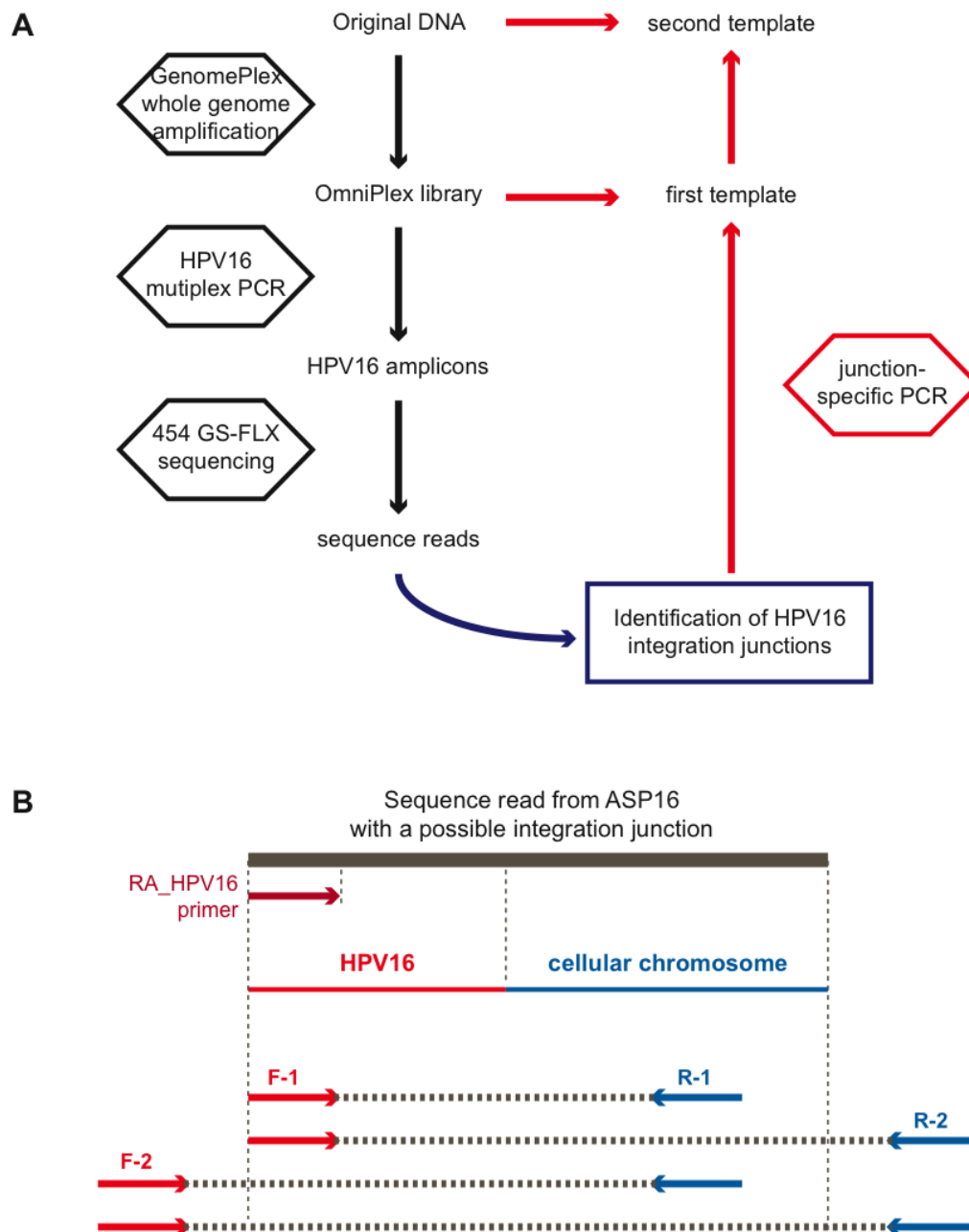


Figure 2.40: Junction-specific PCR for verification of HPV16 integration junctions identified in ASP16 sequence reads. **Panel A** shows the DNA amplification steps for HPV16 amplicon preparation and junction-specific PCR templates. The original DNA is amplified by GenomePlex whole genome amplification kit, to produce the OmniPlex library. HPV16 amplicons for Roche/454 GS-FLX sequencing were prepared from the OmniPlex libraries of each DNA. For the junction-specific PCR, the OmniPlex library was used first as template. After positive PCR results, the original DNA was used as template for PCR, if enough DNA was available. **Panel B** shows the primer design for junction-specific PCR. The top line represents an ASP16 sequence read containing a potential HPV16 integration junction. Two HPV16-specific forward primers, F-1 and F-2, were selected. F-1 binds to HPV16 inside the ASP16 sequence read, and F-2 binds upstream of the ASP16 read. Two reverse primers, R-1 and R-2, were selected in the flanking cellular sequences. R-1 binds the cellular DNA within the ASP16 read, while R-2 binds outside this area. The primers were combined for PCR in four combinations.

Cell line MRI-H186

MRI-H186 cells contain about 30-40 copies of HPV16 DNA integrated into chromosome 8q24 near the c-myc oncogene (Xu, 2010). Two HPV16 3' junctions have been identified, with the major junction at pos. 2754 and the minor junction at pos. 1224. The major junction is present as two integration variants, A and A⁺ (Xu, 2010). Variant A is a truncated HPV16 genome, whereas variant A⁺ includes a complete HPV16 genome (Figure 2.41). An integration percentage of 49% has been determined for this cell line by E2/E6 qPCR.

With the previous knowledge about the HPV16 breakpoint positions, the major and minor integration junctions were expected in the sequences in primer groups E06 (2723) and E19 (1115), respectively. The sequence alignments of primer group E06 in both ASP16-3 and ASP16-4 identified viral-cellular sequences with the major integration junction, together with purely viral sequences (Figure 2.42). The major integration was also detected in sequences of primer E27 (2628), located 95 bp upstream of primer E06. All informative sequences of primer E19, however, did not contain the minor integration junction. The results are summarized in Table 2.19.

Concerning the major 3' junction at pos. 2754, two aspects for junction determination by ASP16 can be noted. First, the distance between primer and breakpoint is short (i.e. 11 nt for sequences in primer group E06), the viral-cellular junction can be determined efficiently in a large number of sequence reads. Second, if the distance between primer and breakpoint is longer (i.e. 104 nt for sequences in primer group E27), a small number of long sequence reads is sufficient to allow an unequivocal determination of a viral-cellular junction. The minor 3' junction at pos. 1223 has been found by PCR experiments as two integration variants, B and A-B (Xu, 2010). In the ASP16 experiments, altogether 118 (24+94) informative sequence reads were obtained from primer E19. However, all reads contained only HPV16 sequences. The reason why ASP16 failed to identify this junction is unclear. However, these results may indicate that variants B and A-B constitute only a minor fraction of the integrated HPV16 genomes in MRI-H186.

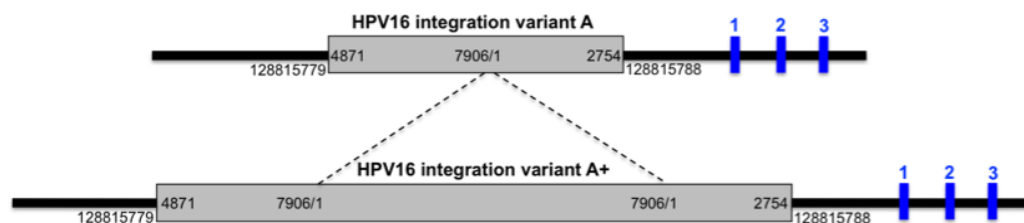


Figure 2.41: Two variants of the major HPV16 integration junctions in MRI-H186. Taken from (Xu, 2010). Variant A and A+ with the major integration junction at HPV16 pos. 2754 are shown. Gray boxes represent HPV16 DNA with positions indicated. Black lines represent cellular DNA. Blue boxes (1, 2, 3) represent the three exons of c-myc gene.

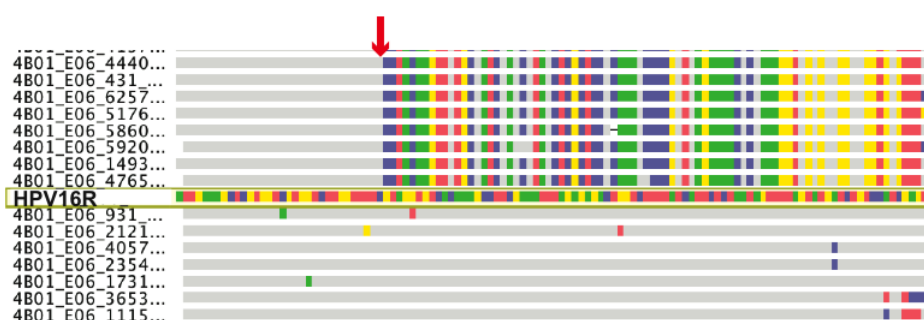


Figure 2.42: Sequence alignment of primer group E06 of MRI-H186 in ASP16-4. The sequences are color-coded. Sequence names are indicated on the left. The reference HPV16R sequence is indicated and all bases are highlighted. For other sequences, bases are highlighted only if they are different from the reference sequence. The sequences above the reference HPV16R contain the integration junction with the HPV16 breakpoint at pos. 2754 (red arrow).

Table 2.19: Features of the major and minor integration junctions of MRI-H186 in ASP16-3 and ASP16-4 sequence reads.

Sample name	MRI-H186						
Integration percentage by E2/E6 qPCR	49%						
Chromosome	8			8			
Cellular breakpoint position (NC_000008.10)	128746606			128675817			
Cellular DNA strand	plus			plus			
HPV16 break point position	(major) pos. 2754			(minor) pos. 1224			
Primer group for junction identification	E06 (2723)		E27 (2628)		E19 (1115)		
Primer positions	2723-2743		2628-2650		1115-1138		
Distance between 3'end of the primer and breakpoint	11 bp		104 bp		86 bp		
Sample ID	2B01*	3B01	4B01	3B01	4B01	3B01	4B01
Junction found?	yes	yes	yes	yes	yes	no	no
# of total reads in this primer group	671	186	485	66	44	283	478
# of reads extending over breakpoint	388	129	444	17	15	24	94
# of reads containing viral-cellular junction	320	112	343	15	14	0	0
# of reads containing HPV16 after breakpoint	68	17	101	2	1	24	94
% of reads containing viral-cellular junction ^(a)	82 %	87 %	77 %	88%	93%	0 %	0 %
Longest sequence read (nt) ^(b)	-	105	130	171	151	-	-

* Experiment ASP16-2 was performed by Bo Xu (Xu, 2010).

(a) The value was calculated from sequence reads extending over the breakpoint. The average value is 85%.

(b) The ≥ 18 nt GPU sequence and the 4-nt barcode were excluded.

Cell line MRI-H196

The HPV16 integration junction in cell line MRI-H196 was identified by RS-PCR (Xu, 2010). The HPV16 breakpoint is located at pos. 3858, integrated into chromosome 11 at pos. 47967861 (NC_000011.9) on plus strand. This leaves the E2 gene (pos.2756-3853) intact in MRI-H196. By E2/E6 qPCR, a low integration percentage of 32% had been determined. In clinical samples, this value would indicate episomal DNA. The E2 primer pair used in E2/E6 qPCR amplifies the segment from pos. 3362 to pos. 3443, which is present in the integrated HPV16 DNA in this cell line.

The sequence reads of primer E16 (3778) should contain the viral-cellular junction sequences. As shown in Table 2.20, all informative sequence reads from this primer group contained the integration junction sequence. In ASP16-3, the longest cellular part was only 8 nt, too short to be identified by Blast without prior knowledge of the integration junction. In ASP16-4, the longest cellular part was 102 nt, long enough to be mapped to chromosome 11.

Table 2.20: Features of the integration junctions of MRI-H196 in ASP16-3 and ASP16-4 sequence reads

Sample name	MRI-H196	
Integration percentage by E2/E6 qPCR	32%	
Chromosome	11	
Cellular breakpoint position (NC_000011.9)	47967861	
Cellular DNA strand	plus	
HPV16 break point position	pos. 3858	
Primer group for junction identification	E16 (3778)	
Primer positions	3778-3798	
Distance between 3'end of the primer and breakpoint	60 bp	
Sample ID	3B02	4B02
Junction found?	yes	yes
# of total reads in this primer group	324	184
# of reads extending over breakpoint	3	36
# of reads containing viral-cellular junction	3 (only 8 nt cellular)	36
# of reads containing HPV16 after breakpoint	0	0
% of reads containing viral-cellular junction ^(a)	100%	100%
Longest sequence read (nt) ^(b)	88	186

(a) The value was calculated from sequence reads extending over the breakpoint.

(b) The ≥ 18 nt GPU sequence and the 4-nt barcode were excluded.

Cell line SiHa

SiHa contains one HPV16 copy integrated into chromosome 13 (Mincheva et al, 1987). The HPV16 3' breakpoint is at pos. 3133 (Meissner, 1999). The sequences of primer E29 (3037) should contain the integration junctions. This was also evident from the pattern of sequence read numbers (Figure 2.38). The results were as expected and are summarized in Table 2.21. Since SiHa cells contain only one integrated HPV16 copy, all informative sequence reads contained the viral-cellular junction. Again, a higher number of long reads in ASP16-4 in comparison to ASP16-3 was observed.

Table 2.21: Features of the integration junction of SiHa in ASP16-3 and ASP16-4 sequences of primer E29 (3037).

Sample name	SiHa	
Chromosome	13	
Cellular breakpoint position (NC_000013.10)	74087558	
Cellular DNA strand	minus	
HPV16 break point position	pos. 3133	
Primer group for junction identification	E29 (3037)	
Primer positions	3037-3058	
Distance between 3'end of the primer and breakpoint	75 bp	
Sample ID	3B03	4B03
Junction found?	yes	yes
# of total reads in this primer group	355	228
# of reads extending over breakpoint	3	16
# of reads containing viral-cellular junction	3	16
# of reads containing HPV16 after breakpoint	0	0
% of reads containing viral-cellular junction ^(a)	100%	100%
Longest sequence read (nt) ^(b)	231	113

(a) The value was calculated from sequence reads extending over the breakpoint.

(b) The ≥ 18 nt GPU sequence and the 4-nt barcode were excluded.

Cell line CaSki

CaSki contains ~600 copies of full-length and rearranged HPV16 DNA (Callahan et al, 1992) integrated as head-to-tail repeats (Baker et al, 1987). More than 10 distinct integration sites have been mapped to different chromosomes (Van Tine et al, 2004). A single transcriptionally active HPV16 integrant has been located on a derivative chromosome 14 in conjunction with chromosome 6 DNA (Van Tine et al, 2004). The viral-cellular junction sequence, found in transcripts and genomic DNA, contains the HPV16 breakpoint at pos. 3728 flanked by chromosome 6 sequences (Jeon & Lambert, 1995; Smits et al, 1991; Van Tine et al, 2004).

Sequence reads of primer E32 (3696) may contain the integration junction. Due to the high number of about 600 HPV16 copies, the chance of finding sequences with the integration junction was estimated to be about 1:600. Among the 773 informative sequences in primer group E32, however, no sequence with integration junction was found. The features for sequence reads of primer E32 are shown in Table 2.22.

Table 2.22: Features of primer E32 sequence reads of CaSki in ASP16-4.

Sample name	CaSki
Chromosome	6
Cellular breakpoint position (NC_000006.11)	45659121
Cellular DNA strand	minus
HPV16 break point position	pos. 3728
Primer group for junction identification	E32 (3696)
Primer positions	3696-3717
Distance between 3'end of the primer and breakpoint	11 bp
Sample ID	4B04
Junction found?	no
# of total reads in this primer group	793
# of reads extending over breakpoint	773
# of reads containing viral-cellular junction	0
# of reads containing HPV16 after breakpoint	3728
% of reads containing viral-cellular junction ^(a)	0%
Longest sequence read (nt) ^(b)	172

(a) The value was calculated from sequence reads extending over the breakpoint.

(b) The ≥ 18 nt GPUTA sequence and the 4-nt barcode were excluded.

CA-07C381

The HPV16 integration junction in this DNA sample had previously been determined in ASP16-2 (Xu, 2010). This sample was re-analyzed in ASP16-4 (ID: 4B05). CA-07C381 contains exclusively integrated DNA (Xu, 2010), see also Table 2.23. In agreement with the known HPV16 breakpoint at pos. 2783, the pattern of sequence read numbers clearly indicated that the breakpoint is situated downstream of primer E06 (2723) (Figure 2.43). All 26 informative E06 sequence reads contained the viral-cellular junction sequence (Table 2.23), consistent with the fully integrated status of HPV16 DNA in this tumor. Comparing ASP16-2 and ASP16-4, the number of informative sequence reads has been significantly increased from 2 to 26.

Table 2.23: Features of primer E06 sequence reads of CA-07C381 in ASP16-4 in comparison with those of ASP16-2

Sample name	CA-07C381	
Integration percentage by E2/E6 qPCR	100%	
Chromosome	12	
Cellular breakpoint position (NC_000012.11)	57686270	
Cellular DNA strand	plus	
HPV16 break point position	pos.2783	
Primer group for junction identification	E06 (2723)	
Primer positions	2723-2742	
Distance between 3'end of the primer and breakpoint	41 bp	
Sample ID	4B05	2B06*
Junction found?	yes	yes
# of total reads in this primer group	42	
# of reads extending over breakpoint	26	2
# of reads containing viral-cellular junction	26	2
# of reads containing HPV16 after breakpoint	0	0
% of reads containing viral-cellular junction ^(a)	100%	100%
Longest sequence read (nt) ^(b)	198	171

* Experiment ASP16-2 was performed by Bo Xu (Xu, 2010).

(a) The value was calculated from sequence reads extending over the breakpoint.

(b) The ≥ 18 nt GPUTA sequence and the 4-nt barcode were excluded.

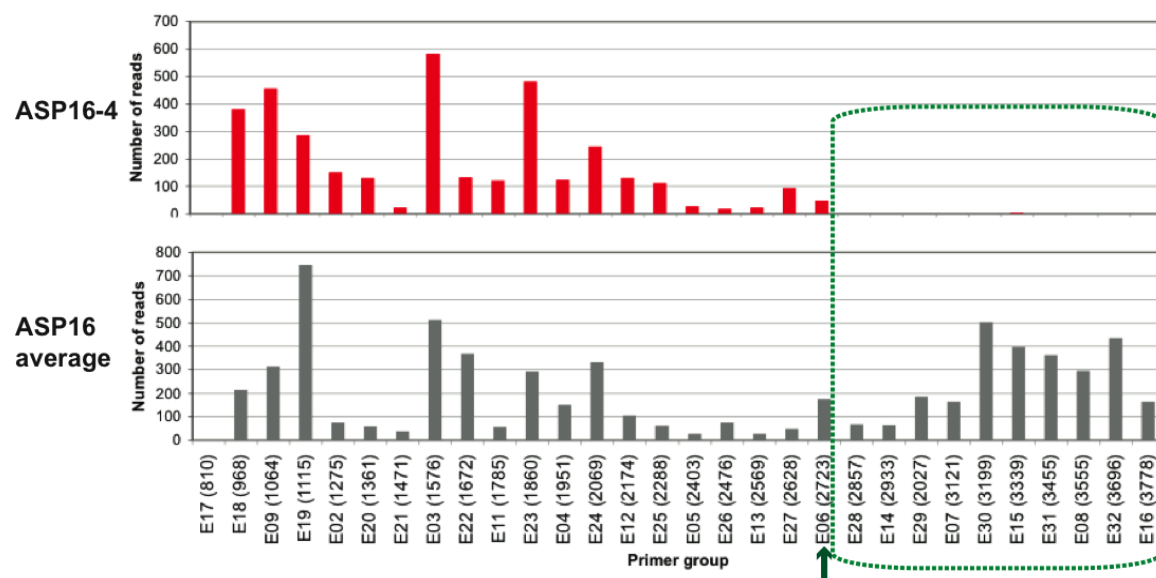


Figure 2.43: Patterns of sequence read numbers in ASP16-4 for CA-07C381. The graph (red) represents the number of reads in ASP16-4 after 28-nt cutoff for CA-07C381. The graph at the bottom (gray) represents the average number of reads in ASP16 after 28-nt cutoff. The data for these plots are shown in Table 2.17. X-axis represents the primer group (with primer name indicated) and Y-axis the number of reads. The possible breakpoint area is circled in green. The green arrow indicates the primer group E06 that contain the HPV16 integration junction sequences with the breakpoint at pos. 2783 (see Table 2.23).

CA-07C368

In ASP16-2, the HPV16 integration junction of sample CA-07C368 was identified in a curious manner (Xu, 2010). Three sequence reads of primer E03 (1576) contained a viral-cellular junction which could not be verified by junction-specific PCR and thus was an artifact. One of the PCR products hybridized strongly with HPV16. Cloning and sequencing revealed a novel viral-cellular junction sequence, which proved to be authentic according to junction-specific PCR. The authentic HPV16 breakpoint is located at pos. 1910, and the HPV16 was found to have integrated into chromosome 1 at pos. 240760698 (NC_000001.10) on plus strand.

CA-07C368 was re-analyzed in ASP16-4 (ID: 4B21). The authentic HPV16 junction was detected among the sequence reads in primer groups E23 (1860) and E11 (1785). The features of these sequence reads are summarized in Table 2.24. The artifact junction sequence of ASP16-2 was not found in the E03 sequence reads of 4B21.

Table 2.24: Features of sequence reads of CA-07C368 in ASP16-4.

Sample name	CA-07C368	
Integration percentage by E2/E6 qPCR	84%	
Chromosome	1	
Cellular breakpoint position (NC_000001.10)	240760698	
Cellular DNA strand	plus	
HPV16 break point position	pos.1910	
Primer group for junction identification	E23 (1860)	E11 (1785)
Primer positions	1860-1879	1785-1805
Distance between 3'end of the primer and breakpoint	32 bp	106 bp
Sample ID	4B21	
Junction found?	yes	yes
# of total reads in this primer group	265	72
# of reads extending over breakpoint	37	1
# of reads containing viral-cellular junction	36	1
# of reads containing HPV16 after breakpoint	1	0
% of reads containing viral-cellular junction ^(a)	97%	100%
Longest sequence read (nt) ^(b)	175	167

(a) The value was calculated from sequence reads extending over the breakpoint.

(b) The ≥ 18 nt GPVA sequence and the 4-nt barcode were excluded.

HSIL-66019 and HSIL-61979

Both HSIL-66019 and HSIL-61979 were isolated from CIN3 lesions and have integration percentage of 100%. In previous experiments, the HPV16 breakpoint of HSIL-66019 was localized between pos. 2451 and 2722 by PCR screening (Steinmeyer, 2009). Because all efforts failed to identify the HPV16 junction by PCR-cloning-sequencing methods (Steinmeyer, 2009), the DNA of HSIL-66019 was analyzed in ASP16-3 (ID: 3B18) and ASP16-4 (ID: 4B06). The possible HPV16 integration junction was detected in sequence reads of 3B18 from primer group E26 (2476) (Figure 2.44, panel A). The junction contained the HPV16 breakpoint at pos. 2516, integrated into chromosome 7 (NC_000007.13) at three possible locations 72454568 (plus strand), 72818793 (plus strand) and 75010338 (minus strand). It was not possible to identify by junction-specific PCR which location is the correct one because these three locations on chromosome 7 contain almost 100% sequence similarity, each location covering a large area of 20-200 kb (Steinmeyer, 2009). The integration junction was confirmed by junction-specific PCR using the original DNA HSIL-66019 as template, and cloning-sequencing of the PCR products (Steinmeyer, 2009). In 4B06, the junction was also detected in sequence reads of primer E26.

For HSIL-61979 (ID: 4B07), an HPV16 integration junction was also detected in primer group E26. This junction sequence was identical to that of HSIL-66019. The original DNA HSIL-61979 was used as template for junction-specific PCR, with the two primer pairs previously used for DNA HSIL-66019 (Steinmeyer, 2009) (Figure 2.44, panel B). Positive products were obtained, cloned and sequenced. The sequences (shown in Appendix A4) were identical to the junction-specific PCR products of DNA HSIL-66019.

It is very unlikely that two unrelated DNA samples have exactly the same HPV16 breakpoint and chromosomal integration site. Review of the sample numbers revealed that both DNA samples were collected from the same woman at a six-month interval, in April 2007 (HSIL-61979) and October 2007 (HSIL-66019). The patterns of sequence read numbers for HSIL-66019 and HSIL-61979 are shown in Figure 2.45. All primer groups from E13 (2569) down to E16 (3778) lacked sequence reads. The primer group E26 (2476), located just upstream of this blank region, contained the viral-cellular junction sequences.

These results clearly demonstrates that the ASP16 strategy is suited both for an appropriate localization of HPV16 integration breakpoints by inspection of sequence read patterns and for exact nucleotide sequence identification of integration junctions. The ASP16 features for the integration junction of samples HSIL-66019 and HSIL-61979 are summarized in Table 2.25.

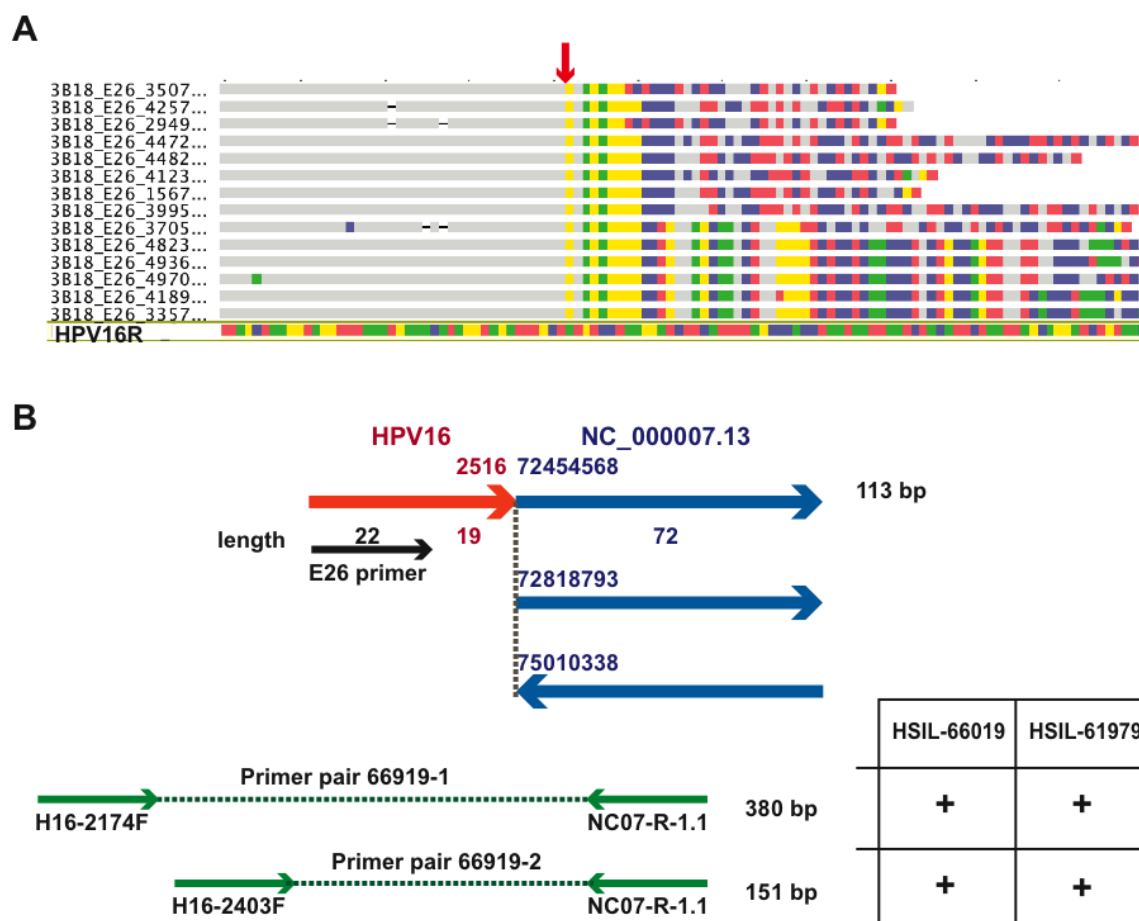


Figure 2.44: HPV16 integration junction in E26 sequences of HSIL-66019 and junction-specific PCR. **Panel A** shows the sequence alignment of E26 sequences of HSIL-66019 (3B18). The sequences are color-coded. Sequence names are indicated on the left. The reference HPV16R sequence is shown at the bottom and all bases are highlighted. For other sequences, bases are highlighted only if they are different from the reference sequence. All 14 reads above the reference HPV16R are composed of HPV16 and cellular sequences, with the HPV16 breakpoint at pos. 2516 (red arrow). **Panel B** shows the junction specific PCR of viral-cellular sequence detected in HSIL-66019 (3B18 and 4B06) and HSIL-61979 (4B07). The composition of the HPV16-cellular sequences is shown on top, with the three possible chromosomal integration sites indicated (blue arrows). The two primer pairs (green arrows) and the expected PCR products are shown below. The PCR results of the original DNAs of both samples are indicated on the right side. The junction-specific PCR for HSIL-66019 had also been performed previously (Steinmeyer, 2009). Primer sequences are given in Materials and Methods.

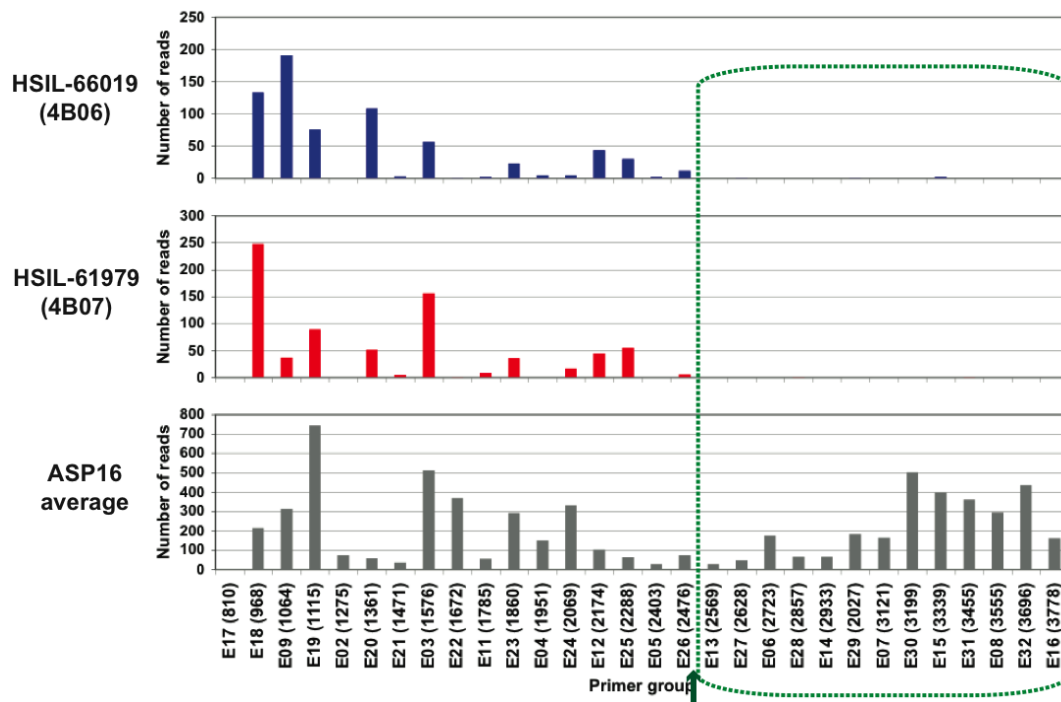


Figure 2.45: Patterns of sequence read numbers in ASP16-4 for HSIL-66019 and HSIL-61979. The graphs represent the number of reads in ASP16-4 after 28-nt cutoff for HSIL-66019 (blue) and HSIL-61979 (red), respectively. The graph at the bottom (gray) represents the average number of reads in ASP16 after 28-nt cutoff. The data for these plots are shown in Table 2.17. X-axis represents the primer group (with primer name indicated) and Y-axis the number of reads. The possible breakpoint area is circled in green. The green arrow indicates primer E26, whose sequences contain the HPV16 integration junction sequences (see Table 2.25).

Table 2.25: Features of the integration junctions of samples HSIL-66019 and HSIL-61979 in ASP16.

Sample name	HSIL-66019		HSIL-61979
Integration percentage by E2/E6 qPCR	100%		100%
Chromosome	7		
Cellular breakpoint position (NC_000007.13)	72454568 or 72818793 or 75010338		
Cellular DNA strand	plus or plus or minus		
HPV16 break point position	pos. 2516		
Primer group for junction identification	E26 (2476)		
Primer positions	2476-2497		
Distance between 3'end of the primer and breakpoint	19 bp		
Sample ID	3B18	4B06	4B07
Junction found?	yes	yes	yes
# of total reads in this primer group	14	11	6
# of reads extending over breakpoint	14	11	6
# of reads containing viral-cellular junction	14	4	6
# of reads containing HPV16 after breakpoint	0	(7) ^(c)	0
% of reads containing viral-cellular junction ^(a)	100%	(36%) ^(c)	100%
Longest sequence read (nt) ^(b)	113	111	117

(a) The value was calculated from sequence reads extending over the breakpoint.

(b) The ≥ 18 nt GPU sequence and the 4-nt barcode were excluded.

(c) Probably due to barcode error.

HSIL-75857

The pattern of sequence read numbers of sample HSIL-75857 (ID: 4B09) (Table 2.17), gave no hint at a possible HPV16 breakpoint location. Therefore, the alignments of all primer groups were studied. A candidate integration junction was detected in the sequence reads of primer E19 (1115) (Figure 2.46, panel A), with HPV16 breakpoint at pos. 1149. The suspected junction was verified by junction-specific PCR with two primer pairs, using the original DNA HSIL-75857 as template (Figure 2.46, panel B). Both PCR reactions were positive, and the products were cloned and sequenced. The sequences (shown in Appendix A4) confirmed the detected HPV16 integration. The features of the sequences of primer E19 of HSIL-75857 are shown in Table 2.26. Although the E2/E6 qPCR value of 100% integration indicated exclusively integrated HPV16 DNA, only 3% of the E16 sequence reads contained the viral-cellular junction, whereas 97% contained purely viral sequences. The reason for this discrepancy is unknown.

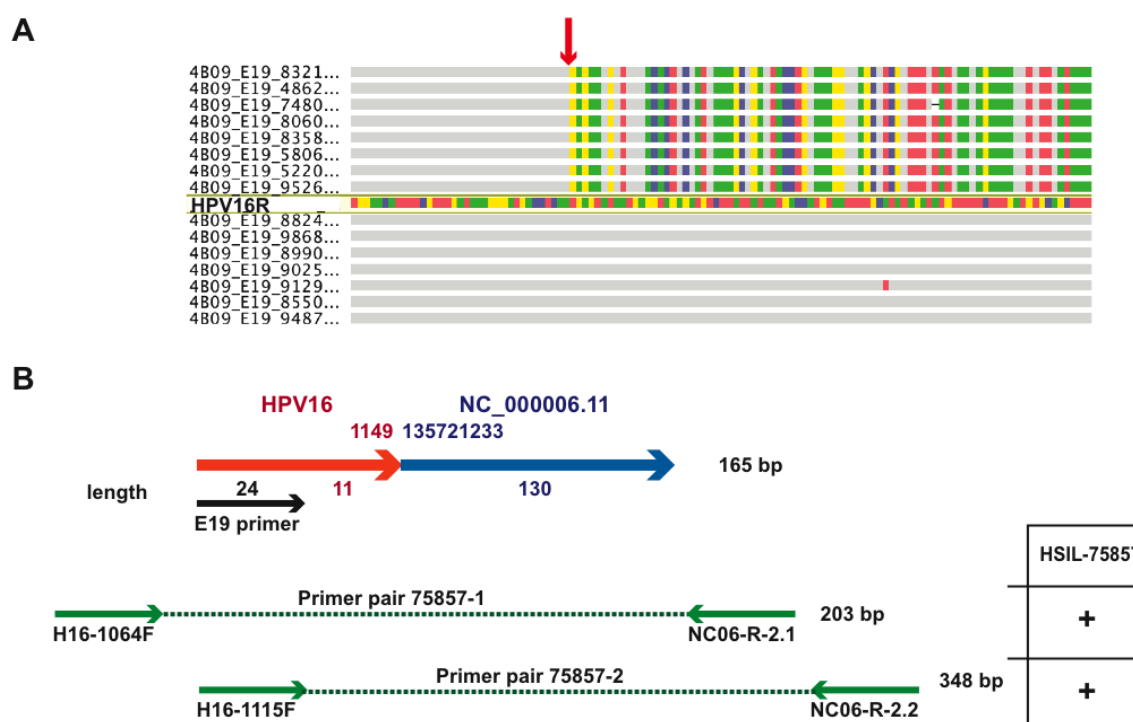


Figure 2.46: HPV16 integration junction in E19 sequences of HSIL-75857 and junction-specific PCR. **Panel A** shows the sequence alignment of E19 sequences of HSIL-75857 (4B09). The sequences are color-coded. Sequence names are indicated on the left. All bases of the reference HPV16R sequence are highlighted. In other sequences, bases are highlighted only if they are different from the reference sequence. All 8 reads above the reference HPV16R contain viral-cellular sequences, with the HPV16 breakpoint at pos. 1149 (red arrow). **Panel B** shows the junction specific PCR of viral-cellular sequence detected in 4B09. The composition of the HPV16-cellular sequences is shown on top. The two primer pairs (green arrows) and the expected PCR products are shown below. The PCR results of the original DNA HSIL-75857 as template are indicated on the right side. Primer sequences are given in Materials and Methods.

Table 2.26: Features of E19 sequence reads of DNA HSIL-75857 in ASP16-4.

Sample name	HSIL-75857
Integration percentage by E2/E6 qPCR	100%
Chromosome	6
Cellular breakpoint position (NC_000006.11)	135721233
Cellular DNA strand	plus
HPV16 break point position	pos. 1149
Primer group for junction identification	E19 (1115)
Primer positions	1115-1138
Distance between 3'end of the primer and breakpoint	11 bp
Sample ID	4B09
Junction found?	yes
# of total reads in this primer group	1453
# of reads extending over breakpoint	1420
# of reads containing viral-cellular junction	41
# of reads containing HPV16 after breakpoint	1379
% of reads containing viral-cellular junction ^(a)	3%
Longest sequence read (nt) ^(b)	165

(a) The value was calculated from sequence reads extending over the breakpoint.

(b) The ≥ 18 nt GPU sequence and the 4-nt barcode were excluded.

CIN2/3-1801

Sample CIN2/3-1801 was analyzed in both ASP16-3 (ID: 3B09) and ASP16-4 (ID: 4B14). The patterns of the sequence read numbers of each primer group in both ASP16-3 and ASP16-4 (Table 2.16 and Table 2.17) gave no hint at any HPV16 breakpoint location. Thus, the sequence alignments of all primer groups were analyzed.

In 3B09, five viral-cellular junction sequences were detected among 11 informative sequence reads of E23 (1860) with HPV16 breakpoint at pos. 1913 (Figure 2.47, panel A). In 4B14, the integration junction sequence was also detected, but only in one sequence read and with only 19 nt of flanking cellular sequence. The features of the detected integration junction in DNA CIN2/3-1801 are summarized in Table 2.27. The possible HPV16 integration junction was examined by junction-specific PCR with two primer pairs using the OminiPlex library and the original DNA as templates. Both templates resulted in positive PCR products for both of the primer pairs (Figure 2.47, panel B). The products from the original DNA template were cloned and sequenced. The sequences (shown in Appendix A4) confirmed the ASP16 sequence data.

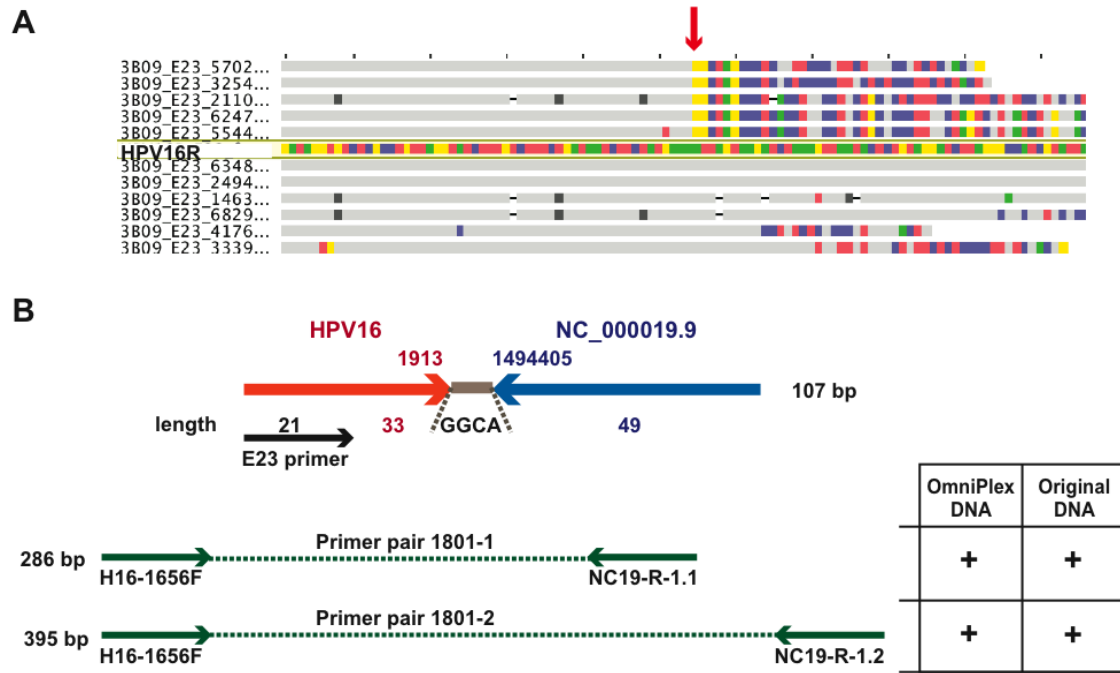


Figure 2.47: HPV16 integration junction in E23 sequences of CIN2/3-1801 from ASP16-3 and junction-specific PCR. Panel A shows the sequence alignment of E23 (1860) sequences of CIN2/3-1801 (3B09). The sequences are color-coded. Sequence names are indicated on the left. All bases of the reference HPV16R sequence are highlighted. In other sequences, bases are highlighted only if they are different from the reference sequence. All reads above the reference HPV16R contain viral-cellular sequences, with the HPV16 breakpoint at pos. 1913 (red arrow). **Panel B** shows the junction specific PCR of viral-cellular sequence detected in 3B09. The composition of the HPV16-cellular sequences is shown on top. The two primer pairs (green arrows) and the expected PCR products are shown below. The PCR results of the OmniPlex library and original DNA of CIN2/3-1801 as templates are indicated on the right side. Primer sequences are given in Materials and Methods.

Table 2.27: Features of the integration junction detected in E23 sequences of CIN2/3-1801 in ASP16.

Sample name	CIN2/3-1801	
Integration percentage by E2/E6 qPCR	32%	
Chromosome	19	
Cellular breakpoint position (NC_000019.9)	1494405	
Cellular DNA strand	minus	
HPV16 break point position	pos.1913	
Primer group for junction identification	E23 (1860)	
Primer positions	1860-1880	
Distance between 3'end of the primer and breakpoint	33 bp	
Sample ID	3B09	4B14
Junction found?	yes	yes
# of total reads in this primer group	252	176
# of reads extending over breakpoint	11	9
# of reads containing viral-cellular junction	5	1 (19 nt celular)
# of reads containing HPV16 after breakpoint	6	8
% of reads containing viral-cellular junction ^(a)	45%	11%
Longest sequence read (nt) ^(b)	107	77

(a) The value was calculated from sequence reads extending over the breakpoint.

(b) The ≥ 18 nt GPU sequence and the 4-nt barcode were excluded.

CIN2/3-4242

The clinical DNA CIN2/3-4242 was derived from a high-grade lesion (integration percentage 79%). The DNA was analyzed in both ASP16-3 (ID: 3B16) and ASP16-4 (ID: 4B20). The pattern of sequence read numbers of 3B16 gave no hints at a possible HPV16 breakpoint. Therefore, the sequence alignments of all primer groups for 3B16 were studied. A possible HPV16 integration junction was detected in one sequence read of primer E15 (3339).

Junction-specific PCRs with four primer pairs were performed using the OmniPlex library and the original DNA as templates (Figure 2.48). From the OmniPlex template, a PCR product was obtained only with primer pair 4242-1, while the original DNA was negative for all primer pairs. The E15 sequence reads of 4B20 did not contain this detected integration junction. The features of the E15 sequence reads of both 3B16 and 3B20 are summarized in Table 2.28. Altogether, the results indicate that the detected viral-cellular sequence was an artifact, and that an HPV16 integration junction for DNA CIN2/3-4242 could not be identified.

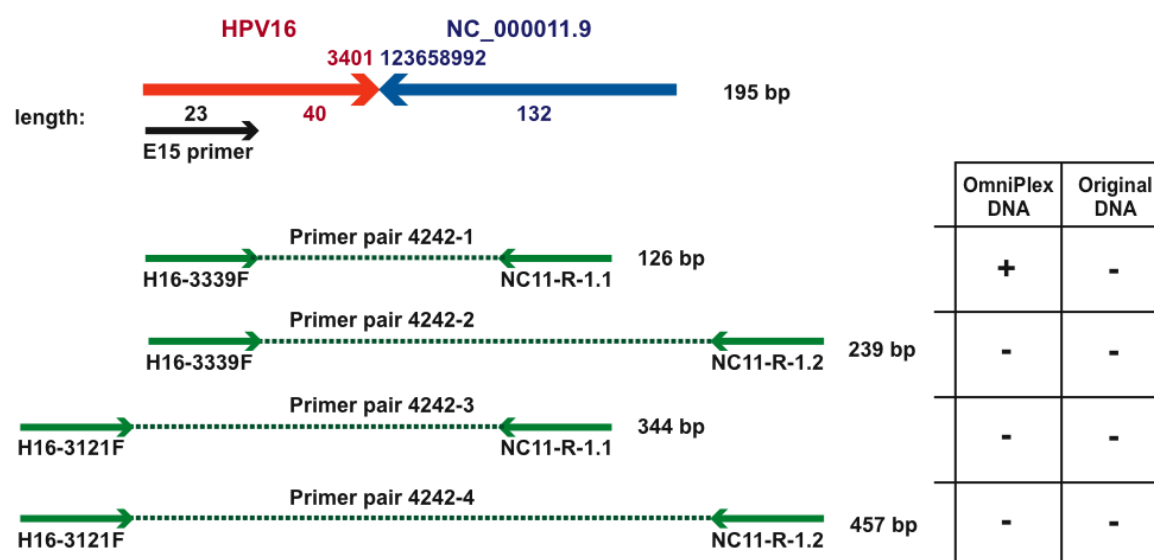


Figure 2.48: Junction-specific PCR of a viral-cellular sequence detected in 3B16. The composition of the HPV16-cellular sequences is shown on top. The four primer pairs (green arrows) and their expected PCR products are shown below. The PCR results of the OmniPlex library and original DNA of CIN2/3-4242 as templates are indicated on the right side. Primer sequences are given in Materials and Methods.

Table 2.28: Features of E15 sequence reads for sample CIN2/3-4242 in ASP16.

Sample name	CIN2/3-4242	
Integration percentage by E2/E6 qPCR	79%	
Chromosome	11	
Cellular breakpoint position (NC_000011.9)	123658992	
Cellular DNA strand	minus	
HPV16 break point position	pos. 3401	
Primer group for junction identification	E15 (3339)	
Primer positions	3339-3361	
Distance between 3'end of the primer and breakpoint	41 bp	
Sample ID	3B09	4B20
Junction found?	false-positive	no
# of total reads in this primer group	197	153
# of reads extending over breakpoint	112	124
# of reads containing viral-cellular junction	(1) ^(a)	0
# of reads containing HPV16 after breakpoint	111	124
% of reads containing viral-cellular junction		0%
Longest sequence read (nt) ^(b)	219	141

(a) False-positive, see text.

(b) The ≥ 18 nt GPUTA sequence and the 4-nt barcode were excluded.

CIN2/3-1503

DNA CIN2/3-1503 was analyzed in ASP16-3 (ID: 3B07) and ASP16-4 (ID: 4B12). The patterns of sequence read numbers gave no hint at a possible HPV16 breakpoint location. Therefore, the sequence alignments of all primer groups were analyzed.

In 3B07, a viral-cellular junction sequence was detected in two sequence reads of primer E19 (1115), with HPV16 breakpoint at pos. 1245 (Figure 2.49, panel A). To examine the authenticity of this junction, junction-specific PCRs with two primer pairs were performed (Figure 2.49, panel B). The OmniPlex library and the original DNA were used as templates. The OmniPlex template gave a positive product only with primer pair 1503-1, while the original DNA CIN2/3-1503 was negative for both primer pairs. The results indicated that the detected viral-cellular sequence is an artifact and not a genuine HPV16 integration junction.

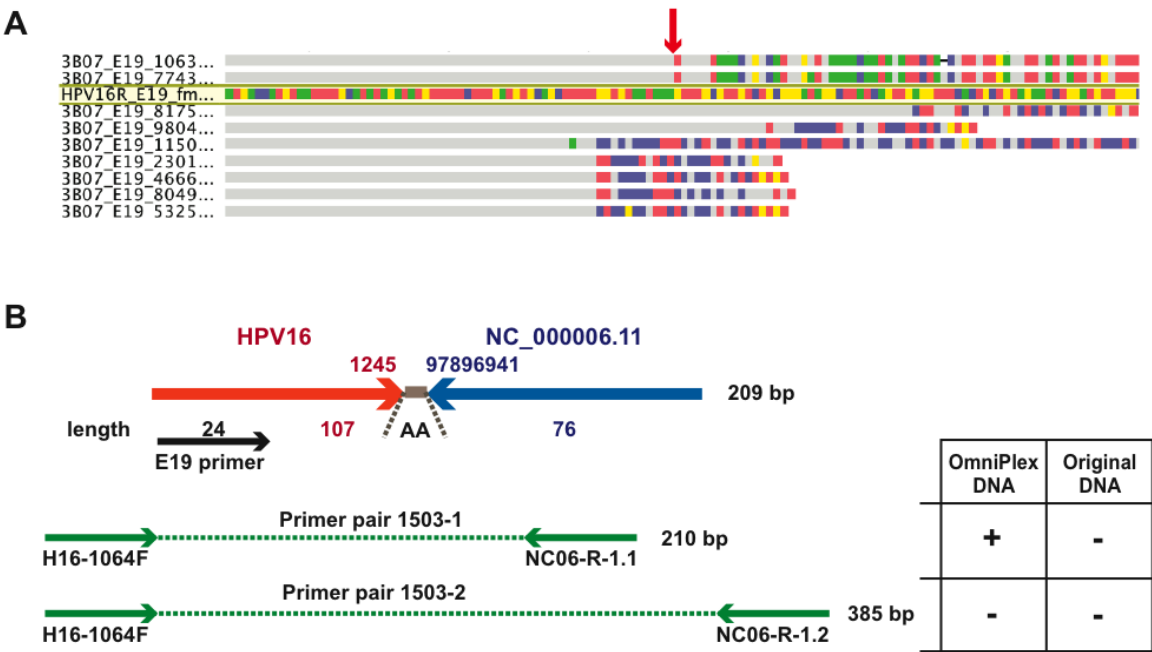


Figure 2.49: Analysis of viral-cellular junction sequences detected in DNA of CIN2/3-1503. Panel A shows the sequence alignment of E19 (1115) sequences of sample 3B07. The sequences are color-coded. Sequence names are indicated on the left. All bases of the reference HPV16R sequence are highlighted. In other sequences, bases are highlighted only if they are different from the reference sequence. The two reads above the reference HPV16R contain viral-cellular sequences, with the HPV16 breakpoint at pos. 1245 (red arrow). **Panel B** shows the junction-specific PCR of viral-cellular sequence detected in 3B07. The composition of the HPV16-cellular sequences is shown on top. The two primer pairs (green arrows) and the expected PCR products are shown below. The PCR results of the OmniPlex library and original DNA of CIN2/3-1503 as templates are indicated on the right side. Primer sequences are given in Materials and Methods.

LSIL- 75022

The DNA LSIL-75022 was derived from a low-grade lesion. Despite its low-grade status, the integration percentage of 100% by E2/E6 qPCR indicates that it should contain fully integrated HPV16 DNA. The viral load was very low (18 copies E7 per 50 ng DNA) compared to other samples (e.g. CA-07C368 with 71923 copies E7 per 50 ng DNA). Consequently, the total number of sequence reads of this DNA obtained in ASP16-4 was very low (see Table 2.17). Based on the pattern of sequence read numbers, possible HPV16 breakpoint positions were localized downstream of primer E23 (1860) and E31 (3455) (Figure 2.50). The sequence alignments of primer groups E23 and E31-E16 were studied. In primer group E23, the alignment shows a number of sequence reads with long CA-rich non-HPV16 sequence (Figure 2.51, panel A). Because the GPU A reverse sequence consists exclusively of C and A residues, it was not possible to figure out whether these sequences are of GPU A or cellular origin. In primer group E31, two out of 28 sequence reads contained a possible viral-cellular sequence of 23 nt (Figure 2.51, panel

B). Blasting these sequences to the human database was unsuccessful in mapping the sequences to a specific chromosome. The sequences of the remaining primer groups were also inspected, but no viral-cellular sequence could be detected. In summary, the breakpoint of the supposed integrated HPV16 DNA in LSIL-75022 could not be identified in experiment ASP16-4.

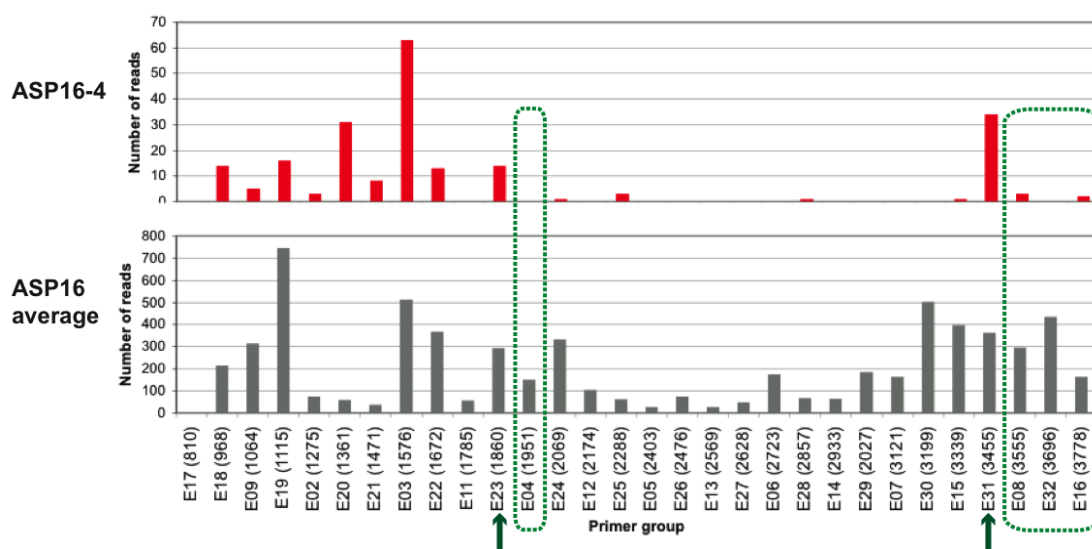


Figure 2.50: Prediction of HPV16 breakpoint locations in DNA LSIL-75022. The graph at the top (red) represents the number of reads for each primer group after 28-nt cutoff of DNA LSIL-75022 (4B08). The graph at the bottom (gray) represents the average number of reads in ASP16. The data for these plots are shown in Table 2.17. X-axis represents the primer group (with primer name indicated) and Y-axis the number of read. The possible breakpoint areas are circled in green. The green arrows indicate primer E23 and E31, whose sequences may contain the HPV16 integration junction sequences.

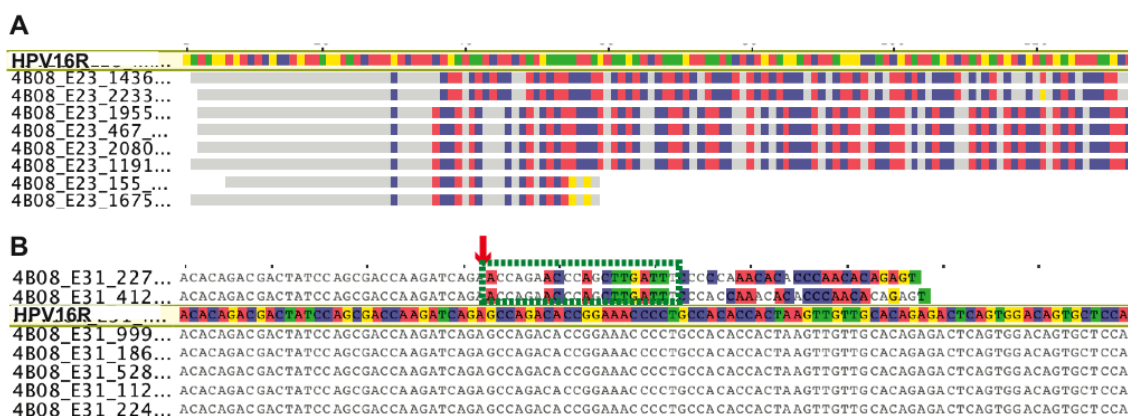


Figure 2.51: Sequence alignment of primer groups E23 and E31 of DNA LSIL-75022. Sequence names are indicated on the left. The sequences are color-coded. All bases of the reference HPV16R sequence are highlighted. In other sequences, bases are highlighted only if they are different from the reference sequence. In **panel A**, all E23 sequences below the reference HPV16R contain long CA-rich non-HPV16 DNA (A and C in red and blue colors). In **panel B**, the sequences above HPV16R contain a possible integration junction with short non-HPV16 sequence (circled in green) with a breakpoint at pos. 3586 (red arrow).

CIN2/3-2219

The DNA CIN2/3-2219 was derived from a high-grade lesion (integration percentage 72%). The patterns of sequence read numbers revealed possible HPV16 breakpoint locations downstream of primers E06 (2723) and E08 (3555) (Figure 2.52). The alignments of sequences from these two primer groups were studied. None of the sequences in these two primer groups was long enough to reach over the next downstream primers (E28 and E32, respectively). No viral-cellular sequence was identified among these sequences. Additionally, the alignments of all other primer groups were also inspected, but no possible integration junction was found.

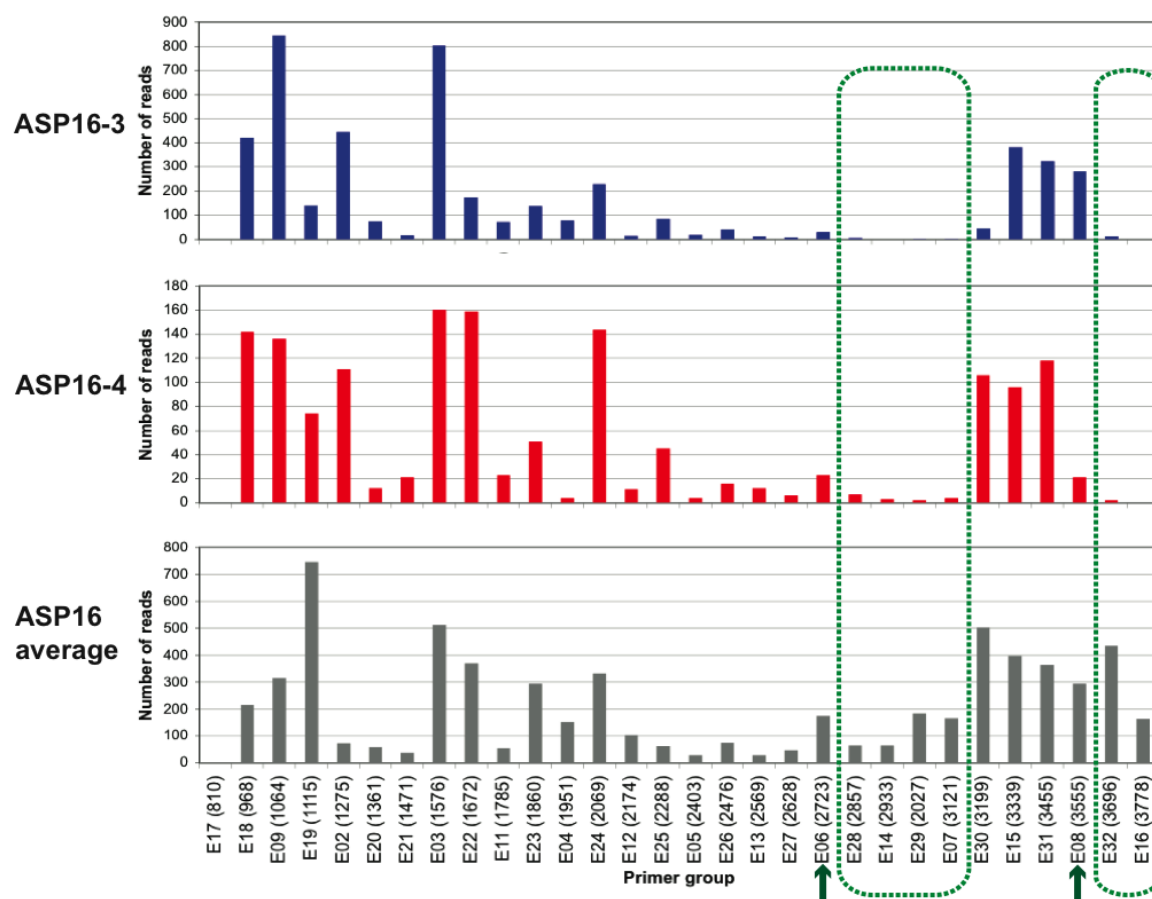


Figure 2.52: Patterns of sequence read numbers for DNA CIN2/3-2219. The graphs represent the number of reads in ASP16-3 (blue) and ASP16-4 (red) after 28-nt cutoff. The graph at the bottom (gray) represents the average number of reads in ASP16-4. The data for these plots are shown in Table 2.16 and Table 2.17. X-axis represents the primer group and Y-axis the number of read. The possible breakpoint areas are circled in green. The green arrows indicate the primer groups (E06 and E08) that may contain sequences with an integration junction.

Clinical DNA samples with no HPV16 breakpoint identified by ASP16

Concerning the remaining eleven DNA samples (Table 2.18), the patterns of sequence read numbers (Table 2.16 and Table 2.17) did not reveal any hint for HPV16 integration breakpoint locations. In addition, after examination of the sequence alignments of all primer groups for these samples, no sequence reads with possible viral-cellular sequences were found. Therefore, no HPV16 integration junctions for these samples could be identified by ASP16.

2.2.2.7 Locations of HPV16 integration junctions

Integration of hr-HPV is an important step toward the development of invasive cancer (Hopman et al, 2004; Schneider-Gadicke & Schwarz, 1986). It usually results in stable expression of the viral oncogenes E6/E7 (Jeon & Lambert, 1995), which is necessary to cause cervical carcinomas (zur Hausen, 1999). Furthermore, it is hypothesized that the integrated HPV DNA can also alter the expression of cellular genes through insertional mutagenesis, such as activation of proto-oncogenes or inactivation of tumor suppressor genes. Therefore, the knowledge of the integration junction locations as well as the cellular sequences directly affected by HPV integration or located in the proximity to the junctions is important to gain more insights into possible insertional mutagenesis effects. For this purpose, the locations of the HPV16 integration junctions and the cellular genes in the vicinity were reviewed by data mining in the NCBI databank. The integration junctions of four clinical DNA samples were newly identified in ASP16-3 and ASP16-4: HSIL-66019, HSIL-61979, HSIL-75857, and CIN2/3-1801. Samples HSIL-66019 and HSIL-61979 were collected from one woman as two independent samples (see previous section) and contain an identical HPV16 integration site, which has already been reviewed for HSIL-66019 (Steinmeyer, 2009). Therefore, in the following, the HPV16 integration junctions in samples HSIL-75857 and CIN2/3-1801 are examined.

For sample HSIL-75857, the HPV16 DNA is integrated into chromosome 6 within an intron of the Abelson helper integration site 1 (AHI1) gene (Figure 2.53). The integration site is located about 98 kb downstream from the 5' end of the first exon of the AHI1 gene, which is transcribed from the opposite strand to the integrated HPV16 DNA. The integration site is located in the proximity of another five cellular genes. The MYB proto-oncogene locates about 181 kb upstream from the integration site, while the

phosphodiesterase 7B (PDE7B) gene is about 451 kb downstream. Another three cellular genes are the HBS1-like (HBS1L), the non-protein coding RNA 271 (NCRNA00271), and the glyceraldehyde-3-phosphate dehydrogenase pseudogene 73 (GAPDHP73).

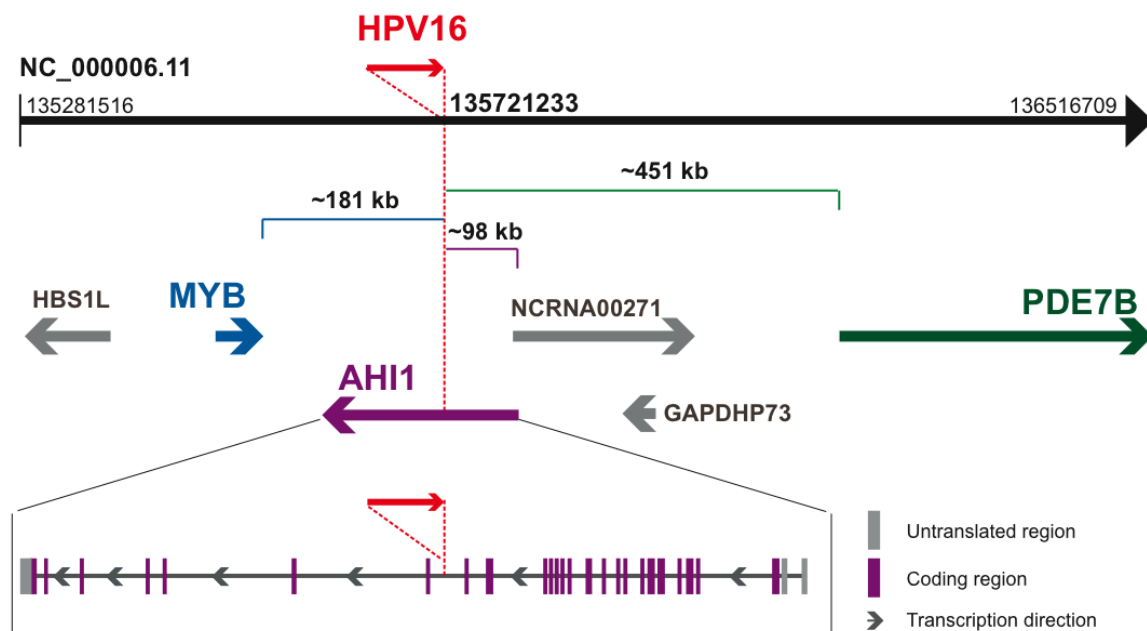


Figure 2.53: Cellular genes in close proximity to the integrated HPV16 of sample HSIL-75857. HPV16 DNA (red arrow) is integrated into chromosome 6 (top black arrow). Gray, blue, purple and green arrows represent coding regions and transcriptional directions of cellular genes, whose names are indicated above the arrows. The inlet at the bottom shows the exon-intron compositions of the AHI1 gene. AHI1: Abelson helper integration site 1. MYB: v-myb myeloblastosis viral oncogene homolog (avian). PDE7B: phosphodiesterase 7B. HBS1L: HBS1-like (*S. cerevisiae*). NCRNA00271: non-protein coding RNA 271. GAPDHP73: glyceraldehyde-3-phosphate dehydrogenase pseudogene 73.

For sample CIN2/3-1801, the HPV16 is integrated in chromosome 19 within an intron of the cellular gene receptor accessory protein 6 (REEP6) (Figure 2.54). The integration site is located 3240 bp downstream from the 5' end of the first exon of REEP6, which is transcribed from the opposite strand to the integrated HPV16 DNA. The proprotein convertase subtilisin/kexin type 4 (PCSK4) gene locates 3998 bp upstream from the integration site, while the adenomatosis polyposis coli 2 (APC2) gene is about 21 kb upstream. There are another three cellular genes in the proximity of the integration site: ADAMTS-like 5 (ADAMTSL5), pseudogene polo-like kinase 5 (PLK5P), and chromosome 19 open reading frame 25 (C19orf25).

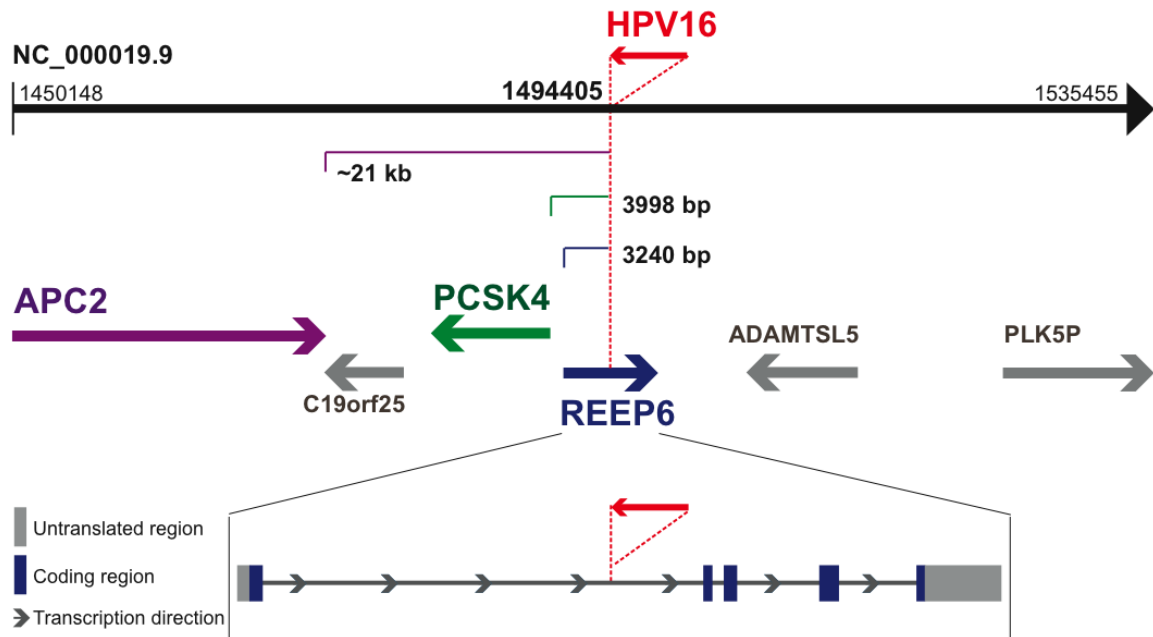


Figure 2.54: Cellular genes in close proximity to the integrated HPV16 of sample CIN2/3-1801. HPV16 DNA (red arrow) is integrated into chromosome 19 (top black arrow). Gray, green and blue arrows represent coding regions and the transcriptional directions of the cellular genes, whose names are indicated above the arrows. The inlet at the bottom shows the exon-intron composition of the REEP6 gene. REEP6: receptor accessory protein 6. ADAMTSL5: ADAMTS-like 5. APC2: adenomatosis polyposis coli 2. C19orf25: chromosome 19 open reading frame 25. PCSK4: proprotein convertase subtilisin/kexin type 4. PLK5P: polo-like kinase 5 (pseudogene).

2.2.2.8 Sequence coverage of the HPV16 E1-E2 region by ASP16 sequence reads

To determine whether the sequence reads from ASP16-3 and ASP16-4 cover the complete analyzed HPV16 E1-E2 area, a “representative E1-E2 sequence” of each DNA sample was created. First, a “representative primer group sequence” was selected from the alignment of each primer group of a DNA sample. Figure 2.55 shows an example. After the representative primer group sequences had been selected for all primer groups of a DNA sample, they were assembled to get the representative E1-E2 sequence of that DNA sample (Figure 2.56). The assembled sequences cover pos. 990-3853, starting with the first nucleotide after primer E18 (pos. 968-989) up to pos. 3853, which is the third nucleotide of the E2 stop codon TGA. Sequences from primer group E17 (810) were not included in the assembled E1-E2 sequences, because they were too short (about 50-60 nt) and never overlapped with sequences of the next primer E18 (968).



Figure 2.55: Example of a representative primer group sequence. The sequence alignment of primer group E08 of sample 3B11 is shown. HPV16R was used as reference sequence (top line). The four bases are highlighted in different color: A in red, C in blue, G in yellow and T in green. A representative primer group sequence (circled in blue) was selected as the sequence containing the most frequent nucleotide at every position, excluding the primer sequence. Red circles mark nucleotides that are different from the most frequent nucleotide at the given position.

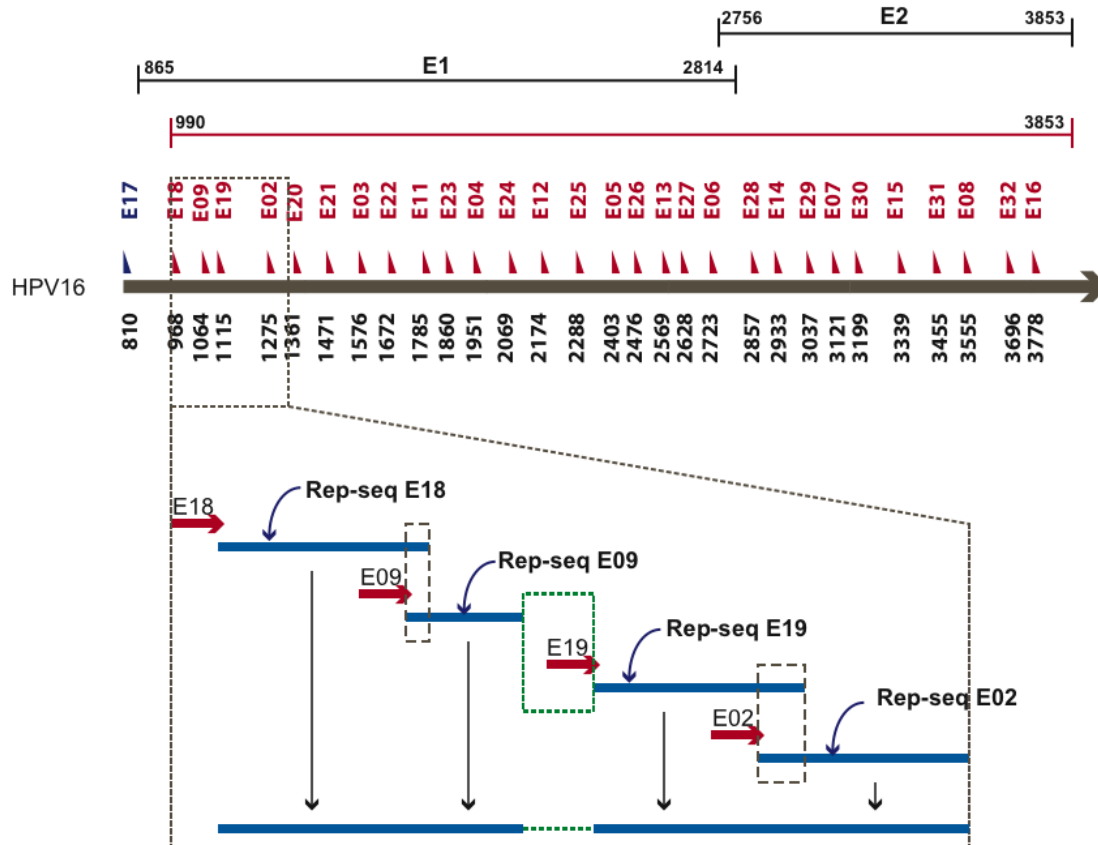


Figure 2.56: Model for assembly of a representative HPV16 E1-E2 sequence for a DNA sample. The locations of all RA_HPV16 primers are indicated along the black arrow. ORFs E1 and E2 are shown above as well as the area of assembled sequences (red line, top). In the lower part, the region from E18 to E02 is enlarged. The arrows represent four primers, for which representative primer group sequences (blue lines, “Rep-seq”) are shown. The “representative E1-E2 sequence” of this DNA sample, assembled from the representative primer group sequences, is shown as blue line at the bottom. The dotted green line indicates a region where no sequence information is available from the ASP16 sequence reads.

From each assembled representative E1-E2 sequence, the sequence coverage was calculated by the following formula:

$$\frac{(\#_of_nucleotides_in_assembled_seq_from_pos._990_to_3853)}{(\#_of_total_nucleotides_from_pos._990_to_3853)} \times 100\%$$

The sequence coverage values for the 25 DNA samples are shown in Figure 2.57. The DNA samples HSIL-66019, HSIL-61979 and LSIL-75022 have very low sequence coverage, compared to others, due to the low number of sequence reads. Excluding these three DNA samples, the average sequence coverage in ASP16-3 and ASP16-4 was 89 %. Compared to the ASP16-2 in which only ~50% sequence coverage could be achieved, the optimized conditions in ASP16-3 and ASP16-4 gave much higher sequence coverage. The distribution of covered and uncovered parts in the assembled representative E1-E2 sequences are illustrated in Figure 2.58. The nucleotide sequences are given in the Appendix A5.

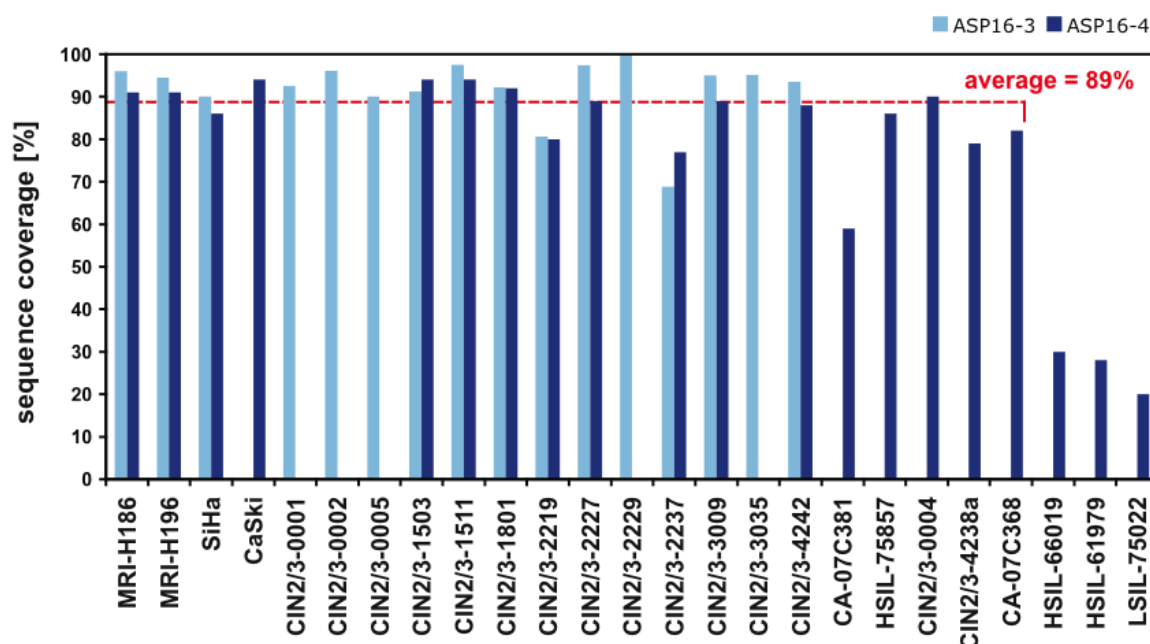


Figure 2.57: Sequence coverage in ASP16-3 and ASP16-4. Each column represents a sequence coverage value (in percent) of a DNA sample in each experiment. The sequence coverage covers the HPV16 E1-E2 area from pos. 990-3853. The columns are grouped by the DNA sample. The average sequence coverage from both experiments is 89%, excluding values from DNA samples HSIL-66019, HSIL-61979 and LSIL-75022.

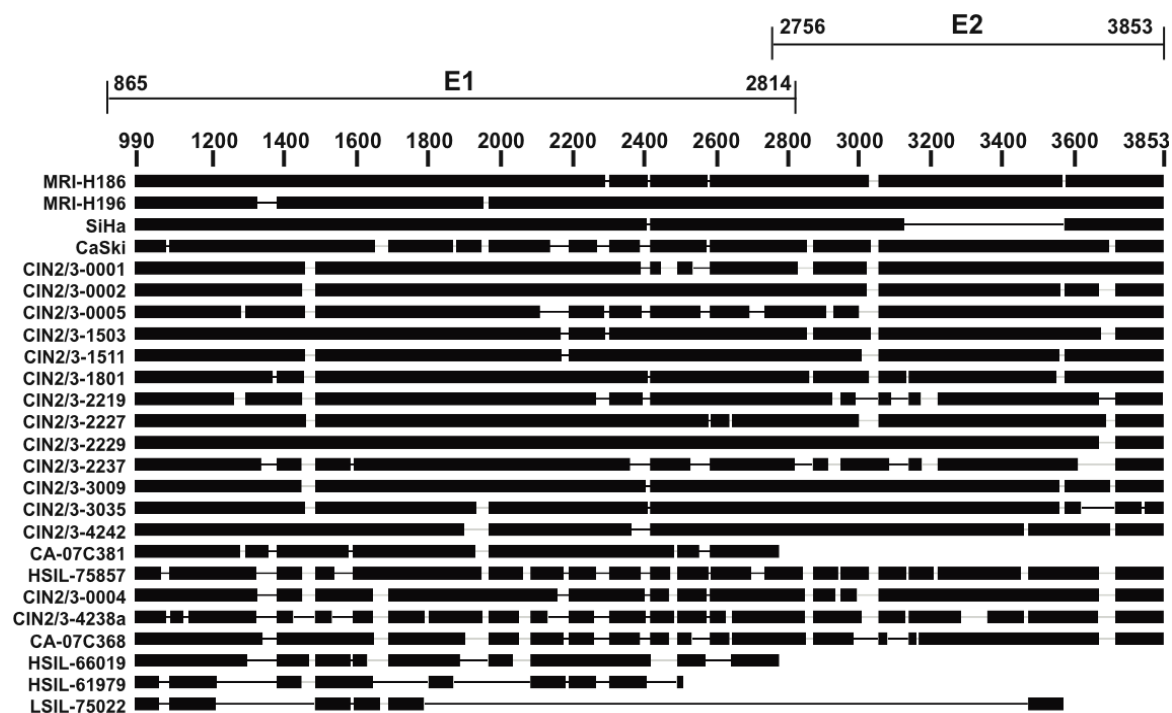


Figure 2.58: Assembled representative E1-E2 sequences of 25 DNA samples in ASP16-3 and ASP16-4. The HPV16R nucleotide positions are indicated on top, the DNA sample names on the left. ORFs E1 and E2 regions are indicated. Each black block represents a covered part for which sequence information could be assembled. Thin lines or blank represent uncovered parts for which no sequence information could be assembled from the ASP16 sequence reads.

2.2.2.9 Nucleotide mutation analysis in HPV16 E1-E2 area and HPV16 variant classification

HPV16 isolates are classified into variants, which are defined as HPVs that are maximally 2% different in their L1 nucleotide sequence (de Villiers et al, 2004). Many HPV16 variants had been identified not only by L1 sequence, but also by other regions such as E6, L2 and URR (Azizi et al, 2008; Casas et al, 1999; Chan et al, 1992; Ho et al, 1993; Kammer et al, 2002; Swan et al, 2005; Yamada et al, 1997; Yamada et al, 1995). Based on geographical distribution, five major phylogenetic branches (also called “lineages”) and one minor branch were identified, designated European (E), Asian (As), African-1 (Af1), African-2 (Af2), Asian-American (AA) and North-American-1 (NA1), respectively (Chan et al, 1992; Ho et al, 1993; Yamada et al, 1995). Linkage of nucleotide polymorphisms between different regions of the same HPV16 isolates had been demonstrated (Casas et al, 1999; Eriksson et al, 1999; Swan et al, 2005). Some HPV16 variants were more often detected in high-grade cervical lesions than others, suggesting an important impact of HPV16 variants in cervical carcinogenesis (Kammer et al, 2000;

Matsumoto et al, 2000; Zehbe et al, 1998). HPV16 variants were shown to have different biochemical and biological properties and these may be one of the reasons for their different carcinogenic risks (Giannoudis & Herrington, 2001).

The assembled representative E1-E2 sequences were used for HPV16 variant classification of the 25 analyzed DNA samples. This was first done for E2. First, a collection of HPV16 variants gathered from the literature were sorted according to their lineage, as shown in Table 2.29, panel A. For each lineage, the nucleotide polymorphism patterns of the collected variants were manually aligned, and the variants were sub-grouped accordingly. As a result, a set of nine “reference nucleotide polymorphism pattern groups” covering the HPV16 E2 ORF (pos. 2756-3853), was created: four groups for the E lineage, two for the AA lineage and one for each As, Af1 and Af2 lineage. A reference nucleotide polymorphism pattern was created for each group, with mandatory and optional nucleotide mutations defined. The mandatory mutations are present in all variants within the same group, while the optional mutations are not. Because no E2 sequence has yet been described for the NA1 lineage, any variant of NA1 lineage could not be identified by this set of reference E2 polymorphism patterns.

Based on the available sequence information in the assembled representative E1-E2 sequences, the twenty-five DNA samples were assigned to the most similar or identical reference pattern, as shown in Table 2.29, panel B. All nucleotide mutations in the E2 ORF region were included in the table. Fifteen samples were successfully assigned to a single pattern, while it was not possible for the remaining 10 samples due to lack of adequate sequence information. Additional thirteen mutation positions, which are not present in the collected HPV16 variant sequences (Table 2.29, panel A), were found in nine DNA samples (Table 2.29, panel B). These mutations may either be variant-specific thus indicating new unknown HPV16 variants, or sample-specific indicating additional, spontaneous mutations occurred in the sample infected with known variants. For DNA samples CIN2/3-3009, CIN2/3-3035 and CIN2/3-4242, their patterns showed the highest similarity to those of AA (1) and AA (2), but do not fit all the mandatory mutations for them to be assigned with either pattern. These three samples may contain HPV16 variants of NA1 lineage because it is closely related to AA lineage (Yamada et al, 1995).

For DNA samples that could not be successfully assigned to a variant/lineage by the E2 sequence, assignment was accomplished by using the HPV16 E6 region. The nucleotide sequence of the NA1 lineage for the E6 region is available (Yamada et al, 1995), therefore all six lineages could now be assigned. Similar to E2, a set of reference nucleotide polymorphism patterns covering the E6 ORF (pos. 83 -559) was created from a collection of 35 HPV16 variants, as shown in Table 2.30, panel A. Out of these 35 variants, thirty-two had been used previously for the E2 reference patterns. Seven reference patterns were created for all six lineages. The HPV16 variants of the E lineage were divided into two pattern groups, Ep and E-T350G, whereas a single reference pattern was assigned to each of the other five lineages (Table 2.30, panel A). The sequences of HPV16 E6 ORF in the 25 DNA samples were determined by PCR amplification using the primers H16-44F and H16-691R (sequences shown in Materials and Methods), cloning and sequencing. The HPV16 E6 sequences of the DNA samples (shown in Appendix A6) were compared to the reference patterns, and the most similar or identical pattern was assigned to all DNA samples (Table 2.30, panel B).

All DNA samples could now be successfully assigned. As assumed from the E2 sequences, the DNA samples CIN2/3-3009, CIN2/3-3035 and CIN2/3-4242 were assigned to NA1 lineage by E6 sequences. No E2 sequence information had been obtained for samples HSIL-66019, HSIL-61979 and CA-07C381 because the breakpoints of the fully integrated HPV16 genomes are located upstream of E2. The E6 sequences showed that samples HSIL-66019 and HSIL-61979 belong to the Ep variant group, and sample CA-07C381 the E-T350G variant group.

The nucleotide polymorphisms in the E1 gene were obtained from the assembled representative HPV16 E1-E2 sequences and are shown in parallel with those of E2 and E6 genes (Table 2.31). Despite incomplete sequence information in the E1-E2 region, nucleotide co-variation in E1, E2 and E6 genes among the same lineage became apparent, especially for lineages NA1 and Af1. Though incomplete, these data provide for the first time the clear linkage of nucleotide variations in E1, E2 and E6 genes of the NA1 lineage.

The effects of the nucleotide variations on the amino acid sequences were investigated. The E6 protein sequences were translated from the complete E6 nucleotide sequences. For E1 and E2 proteins, the assembled nucleotide representative E1-E2 sequences were

translated into amino acid residues. Because the gaps in the assembled sequences prevent correct-frame translation, “N” nucleotides were added to these gaps to make continuous sequences. The amino acid changes of E6, E1 and E2 protein of all DNA samples are shown in Table 2.32. No mutation inducing premature stop codon (nonsense mutation) was observed. The most frequently detected non-synonymous mutations were E6 L83V (in 16 of 25 samples), E1 S220T (in 12 of 25 samples) and E2 P219S (in 14 of 20 samples). The number of non-synonymous mutations is similar to the number of silent mutations in most cases (European lineage ORFs E6, E1 and E2; NA1/Af1/Af2 lineages ORFs E6 and E1). In ORF E2 of NA1/Af1/Af2 lineages, however, the number of non-synonymous mutations is 3-4 fold higher than that of silent mutations.

Table 2.29: Nucleotide polymorphisms in HPV16 ORF E2 and variant assignment.

A: Reference patterns of nucleotide polymorphisms in ORF E2.

[illegible]

Table 2.30: Nucleotide polymorphisms in HPV16 ORF E6 and variant assignment.

A: Reference patterns of nucleotide polymorphisms in ORF E6.

Pattern group	Pattern group by E2 ORF	Original isolate/variant name from the source	Source	8	1	1	1	1	1	1	1	1	2	2	2	3	3	3	4	4	4	5	5		
(a)	(b)		(b)	3	0	3	3	4	4	7	7	8	5	8	8	3	5	6	0	3	5	0	3		
				9	1	2	3	5	2	6	8	3	6	6	9	5	0	2	3	3	7	5	2		
HPV16R					A	T	A	G	C	G	A	G	T	T	C	T	A	C	T	A	A	G	T	T	
Ep	E (1)	Ep	Ref 7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-		
Ep	E (1)	isolate Qv17722E (accession: AY686584)	Ref 3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	-	-		
Ep	E (2)	Ep-c	Ref 7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-		
Ep	E (2)	H3	Ref 9	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
Ep	E (2)	K1	Ref 9	-	-	-	-	#	C	A	-	-	-	-	-	-	-	-	-	-	-	-	-		
Ep	E (3)	K6	Ref 9	-	-	-	-	#	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-		
Ep	E (3)	Ep-a	Ref 7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-		
Ep	E (3)	isolate Qv15521E (accession: AY686581)	Ref 3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-		
Ep	E (4)	Ep-b	Ref 7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-		
Ep	E (4)	T1	Ref 9	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	-		
E-T350G	E (2)	E-G350	Ref 7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (2)	clone 114/K (accession: EU118173)	Ref 5	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (2)	isolate Qv16936E (accession: AY686580)	Ref 3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (2)	H2	Ref 9	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (2)	C1	Ref 9	-	-	-	-	#	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (2)	E-C109G	Ref 7	-	C	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (2)	isolate Qv18158E (accession: AY686583)	Ref 3	-	-	-	-	-	-	A	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (2)	isolate 16W12E (accession: AF125673)	Ref 3	-	-	-	-	-	-	-	-	-	T	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (2)	K2	Ref 9	-	-	-	-	#	C	A	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (2)	M2	Ref 9	-	-	-	-	#	C	A	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (4)	E-G131G	Ref 7	-	-	G	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (4)	CaSki cell line (accession: U89348)	Ref 6	-	-	G	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
E-T350G	E (4)	isolate European German (AF536179)	Ref 3	-	-	G	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-		
As	As	As	Ref 7	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	-	-	-	-		
As	As	S1	Ref 9	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	-	-	-	-		
AF1	AF1	AF1	Ref 7	C	-	-	C	G	T	-	-	-	-	-	A	G	T	-	-	-	-	-	-		
AF1	AF1	AF1 (accession: AF472508)	Ref 3	C	-	-	C	G	T	-	-	-	-	-	A	G	T	-	-	-	-	-	-		
AF1	-	IS.398	Ref 12	-	-	G	-	G	T	-	-	-	-	-	A	G	T	G	-	-	-	-	-		
AF1	-	IS.347	Ref 12	-	-	-	-	G	T	-	-	-	-	-	A	G	T	G	-	-	-	-	-		
AF2	AF2	AF2 (accession: AF472509)	Ref 3	-	C	-	T	G	T	-	-	-	-	-	A	G	T	-	-	G	-	-	-		
AF2	AF2	AF2	Ref 7	-	C	-	T	G	T	-	-	-	-	-	A	G	T	-	-	G	-	-	-		
NA1	-	OR.3136	Ref 11	-	-	-	-	-	T	-	-	-	-	-	A	G	T	G	-	-	-	-	-		
AA	AA (1)	AA	Ref 7	-	-	-	-	-	T	-	-	-	-	-	A	G	T	G	-	-	-	-	G		
AA	AA (1)	AA (accession: AF402678)	Ref 3	-	-	-	-	-	T	-	-	-	-	-	A	G	T	G	-	-	A	-	G		
AA	AA (2)	K4	Ref 9	-	-	-	-	-	T	-	-	-	G	-	A	G	T	G	-	-	-	C	G		
Reference patterns for E6 polymorphisms (c)				Ep	-	-	-	-	-	C	A	-	-	-	-	-	-	G	-	-	-	C	C	-	
				E-T350G	-	C	G	-	-	-	C	A	-	-	T	-	-	-	-	G	-	-	-	-	-
				As	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	-	-	-	-	-	-
				AF1	C	-	G	C	G	T	-	-	-	-	-	-	-	A	G	T	G	-	-	-	-
				AF2	-	C	-	T	G	T	-	-	-	-	-	-	-	A	G	T	-	-	G	-	-
				NA1	-	-	-	-	T	-	-	-	-	-	A	G	T	G	-	-	-	-			
				AA	-	-	-	-	T	-	-	-	G	-	A	G	T	G	-	-	A	-	C		

B: Nucleotide polymorphisms in ORF E6 of 25 DNA samples in comparison with the reference patterns.

DNA sample	Nucleotide pos. (d)	8	9	1	1	1	1	1	1	1	1	2	2	2	2	3	3	3	4	4	4	4	5	5	Patterns assigned by E2 (\$)	
	E6 amino acid residue (f)	8	9	1	1	1	1	2	2	2	2	5	6	6	7	8	8	9	1	1	1	1	1	1		
	E6 aa change (e)	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*	*		
	HPV16R	A	A	T	A	G	C	G	A	G	T	T	C	T	A	T	C	T	A	C	A	G	A	T		T
Reference patterns (c)	Ep	-	-	-	-	-	-	-	C	A	-	-	-	-	-	-	-	G	-	-	-	-	-	C	C	-
	E-T350G	-	-	C	G	-	-	-	C	A	-	-	T	-	-	-	-	G	-	-	-	-	-	-	-	-
	As	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	Af1	C	-	-	G	C	G	T	-	-	-	-	-	A	G	-	T	G	-	-	-	-	-	-	-	-
	Af2	-	-	C	-	T	G	T	-	-	-	-	-	A	G	-	T	G	-	-	G	-	-	-	-	-
	NA1	-	-	-	-	-	-	T	-	-	-	-	-	-	A	G	-	T	G	-	-	-	-	-	-	-
	AA	-	-	-	-	-	T	-	-	-	-	G	-	A	G	-	T	G	-	-	A	-	-	C	G	
CIN2/3-0004	Ep	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	E (1)
CA-07C368	Ep	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	E (1)
MRI-H186	Ep	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	E (3)
CIN2/3-2229	Ep	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	E (4)
CIN2/3-4238a	Ep	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	E (4)
HSIL-66019	Ep	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	not possible
HSIL-61979	Ep	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	not possible
CIN2/3-0001	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	E (2)
CIN2/3-0002	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	E (2)
CIN2/3-1503	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	E (2)
CIN2/3-1801	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	E (2)
CIN2/3-2227	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	E (2)
HSIL-75857	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	T	-	-	-	-	-	E (2)
MRI-H196	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	G	-	-	-	-	-	-	-	E (2)
SiHa	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	C	-	-	-	E (1 or 2)
CaSki	E-T350G	-	-	G	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	E (4)
CA-07C381	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	not possible
LSIL-75022	E-T350G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	-	-	-	-	-	-	-	E (1 - 4)
CIN2/3-3009	NA1	-	T	-	-	-	T	-	-	-	-	-	A	G	-	T	G	-	-	-	-	-	-	-	-	(similar to AA)
CIN2/3-3035	NA1	-	-	-	-	-	T	-	-	-	-	-	A	G	-	T	G	-	-	-	-	-	-	-	-	(similar to AA)
CIN2/3-4242	NA1	-	-	-	-	-	T	-	-	-	-	-	A	G	-	T	G	-	-	-	-	-	-	-	-	(similar to AA)
CIN2/3-0005	Af1	-	-	-	G	-	G	T	-	-	-	-	A	G	G	T	G	-	-	-	-	-	-	-	-	Af1 or Af2
CIN2/3-1511	Af1	-	-	-	G	-	G	T	-	-	-	-	A	G	G	T	G	-	-	-	-	-	-	-	-	Af1 or Af2
CIN2/3-2219	Af1	C	-	-	-	C	G	T	-	-	-	-	A	G	-	T	-	-	-	-	-	-	-	-	-	Af1
CIN2/3-2237	Af2	-	-	C	-	T	G	T	-	-	-	-	A	G	-	T	-	-	G	-	-	-	-	-	-	Af2

(a)-(e) See footnotes of Table 2.29. (b) additional references: Ref 11 (Yamada et al, 1995), and Ref 12 (Eriksson et al, 1999). (f) Because the early promoter is located at nucleotide 97, ORF E6 is translated from the second ATG codon (nucleotide 104). (§) Referred to Table 2.29.

Table 2.32: Amino acid residue variations in HPV16 E6, E1 and E2 proteins in 25 DNA samples.

[illegible]

(§) The positions of the silent mutations are given in Table 2.29 panel B, Table 2.30 panel B and Table 2.31.

(§§) Referred to Table 2.30. “#” indicates no sequence information.

2.2.2.10 Summary of ASP16 results for HPV16 integration and mutation analysis

The integration junctions identified in ASP16-3 and ASP16-4, the sequence coverage and the assigned variants/lineages of each sample are shown in Table 2.33. All four cell lines belong to the European lineage. The 21 clinical DNA samples were obtained from 20 women, with samples HSIL-66019 and HSIL-61979 collected from the same woman. Among these 20 individuals, 13 contain HPV16 variants of European lineage, 4 of African and 3 of North-American.

Table 2.33: Summary of the integration and sequence analyses of ASP16-3 and ASP16-4.

Sample name ^(a)	Integration %	HPV16 breakpoint	Chromosome - position - strand - NC.version	HPV16 variant ^(b)	Sequence coverage
<i>Cell lines</i>					
MRI-H186	49	2754	8 - 128746606 - plus - NC_000008.10	Ep	96%
MRI-H196	32	3858	11 - 47967861 - plus - NC_000011.9	E-T350G	95%
SiHa		3133	13 - 74087558 - minus - NC_000013.10	E-T350G	90%
CaSki				E-T350G	94%
<i>Clinical samples for which HPV16 integration breakpoints were identified by ASP16 *</i>					
CA-07C381	100	2783	12 - 57686270 - plus - NC000012.11	E-T350G	59%
CA-07C368	84	1910	1 - 240760698 - plus - NC000001.10	Ep	82%
HSIL-66019	100	2516	7 - 72454568 - plus - NC_000007.13, or 7 - 72818793 - plus - NC_000007.13, or 7 - 75010338 - minus - NC_000007.13	Ep	30%
HSIL-61979	100				28%
HSIL-75857	100	1149	6 - 135721233 - plus - NC_000006.11	E-T350G	86%
CIN2/3-1801	72	1913	19 - 1494405 - minus - NC_000019.9	E-T350G	92%
<i>Clinical samples for which HPV16 integration breakpoints were not identified by ASP16 **</i>					
CIN2/3-0004	82			Ep	90%
CIN2/3-2229	76			Ep	100%
CIN2/3-4238a	64			Ep	79%
CIN2/3-0001	76			E-T350G	93%
CIN2/3-0002	74			E-T350G	96%
CIN2/3-1503	74		false-positive	E-T350G	94%
CIN2/3-2227	99			E-T350G	97%
LSIL-75022	100			E-T350G	20%
CIN2/3-3009	81			NA1	95%
CIN2/3-3035	74			NA1	95%
CIN2/3-4242	79		false-positive	NA1	94%
CIN2/3-0005	95			Af1	90%
CIN2/3-1511	87			Af1	98%
CIN2/3-2219	72			Af1	81%
CIN2/3-2237	98			Af2	77%

(a) The sample names include the histological/cytological status as a prefix. CA: cervical carcinoma CIN: cervical intraepithelial neoplasia. LSIL/HSIL: low/high-grade squamous intraepithelial lesion.

(b) The variant assignment is based on E6 polymorphisms. See also Table 2.30.

* Sorted by histological/cytological status.

** Sorted by lineage/variants.

3. Discussion

3.1 Complete sequence of integrated HPV68b in ME180 and ME180R

The complete sequences of the integrated HPV68b DNA in cell lines ME180 and ME180R have been determined in this study. In comparison to the previous sequence of cloned restriction fragment AA13.1 (Reuter et al, 1991), the newly determined HPV68b(int)-ME180 is about 6.1 kb larger. The size difference corresponds to the probably cloning-induced rearrangement of AA13.1 where almost the complete 5'copy was deleted (Figure 2.3). The new sequence HPV68b(int)-ME180 contains two incomplete HPV68b copies, integrated in head-to-head arrangement. In cervical carcinoma cell lines, HPVs were found to be integrated as a single incomplete copy such as in SiHa (Meissner, 1999), or as a disrupted copy flanking an additional complete copy as in MRI-H186 (Xu, 2010) (see Figure 2.41) or flanking multiple complete copies as in CaSki (Baker et al, 1987; Callahan et al, 1992). These multiple integrated HPV copies were found in tandem head-to-tail arrangement. To our knowledge, the integrated HPV68b in ME180 is the only known case of head-to-head arrangement of integrated HPV in cervical carcinoma cell lines. Because the URR, and ORFs E6 and E7 in both 5'copy and 3'copy are intact, the viral oncogenes E6/E7 can be expressed from both copies.

In the sub-line ME180R which had been selected for resistance to growth-inhibition by TNFalpha (Pfreundschuh et al, 1989), the sequence HPV68b(int)-ME180R contains two large deletions, compared to ME180 (Figure 2.6). As the result, only URR, and ORFs E6 and E7 in the 5'copy are intact. It was investigated in this study whether alterations of the integrated HPV68b DNA are a recurrent phenomenon associated with selection of ME180 cells for resistance to TNFalpha. Two new TNFalpha-resistant variants, ME180-2A and ME180-3A, were selected which showed no alteration in the integrated HPV68b DNA. Thus, it was concluded that the two large deletions in ME180R are not necessary for TNFalpha resistance phenotype. Other cellular variations may contribute to this phenotype (Manchester et al, 1993; Nishikawa et al, 1992).

3.2 Full-length HPV68b genomes

In this study, a full-length (complete) HPV68b genome, HPV68b-CIN2, and a partially deleted genome, HPV68b-CIN2-Del, have been isolated from a CIN2 DNA sample. Compared to the full-length genome HPV68b-CIN2, the HPV68b-CIN2-Del genome carries a 1229-bp deletion in ORF E1, resulting in a premature stop codon of E1 immediately downstream of the deletion site (Figure 2.15). Since functional E1 protein is necessary for viral episomal genome maintenance but unnecessary for integrated HPV, it was hypothesized that the HPV68b-CIN2-Del genome was likely integrated rather than episomal. This is supported by genomic Southern blot analysis, which indicated a probably integrated status. Because no additional CIN2 genomic DNA for further studies was available, it was impossible to proof this assumption and to determine whether the full-length HPV68b-CIN2 genome is integrated or episomal. Furthermore, it remained unclear whether the two genomes were present in the same cell or originated from different cells of the CIN2 sample.

The presence of intact genome and E1-disrupted genome of HPV62 in a single cervical sample of normal cytology had been reported (Fu et al, 2004), where a single nucleotide in the 5' region of E1 was deleted in the disrupted genome resulting in a premature stop codon. Single nucleotide deletions in E1 ORF that caused premature stop codons were also detected in the reference genome sequences of HPV16, HPV53, HPV56, and HPV72 (Fu et al, 2004). Although the HPV68b-CIN2-Del carries a much larger deletion in ORF E1, it resulted in the same consequence of a premature stop codon. Altogether, these data show that the disruption in ORF E1 maybe rather common among genital HPVs.

The nucleotide sequence of HPV68b-CIN2 shows 99% similarity to HPV68b-ME180, 93% similarity to HPV68a subtype (accession DQ080079), and 99% similarity to the recently published complete nucleotide sequence of the HPV68b variant HPV68b-EU918769 (Wu et al, 2009). The HPV68b-CIN2 and HPV68b-EU918769 genomes are the only two HPV68b variants whose complete nucleotide sequences are available up to now.

3.3 HPV68 subtypes and variants

In an effort to determine HPV68 subtypes/variants among the eleven HPV68-positive clinical samples analyzed in this study, the nucleotide sequences (491 bp) of URR have been determined and compared with other published sequences (Table 2.10). The phylogenetic tree, constructed from these URR sequences, including HPV68b-CIN2, HPV68b-ME180, HPV68b-EU918769 and other HPV68b variants, clearly shows two deep dichotomic branches for HPV68a and HPV68b subtypes (Figure 2.21).

Sample Reims-12 contains subtype HPV68a. The analyzed URR sequence of sample Reims-12 is identical to that of HPV68a-DQ080079. Since both HPV68a genomes were isolated from cervical lesions of French women, it is likely that these two HPV68a genomes belong to the same HPV68a variant present in France. DQ080079 is the only published HPV68a sequence so far (Longuet et al, 1996).

The 41 samples in the study of (Calleja-Macias et al, 2005) had been collected from Brazil, Mexico, the United States, South Africa, Hong Kong and Scotland, while all eleven samples analyzed in this study were collected in France. All 41 samples (100%) in the study of (Calleja-Macias et al, 2005) contain subtype HPV68b, whereas 10 out of 11 samples (91%) analyzed in this study contain subtype HPV68b. It is clear from these observations that subtype HPV68b is more widely distributed and more frequently found than subtype HPV68a. HPV68a may be more prevalent in France than in other countries. It will be interesting to analyze whether HPV68a and HPV68b differ in their propensity to establish persistent infections, as has been shown for HPV16 variants (Schiffman et al, 2010).

3.4 Optimized ASP16 strategy and computer analysis programs

The ASP16 strategy was developed for massively parallel sequencing of HPV16 DNA in clinical samples with the main purpose of determining HPV16 integration junctions (Xu, 2010). It employs one of the next generation sequencing technologies, Roche/454 GS-FLX pyrosequencing (<http://www.454.com/>).

Because of the high output volumes for each sequencing run by Roche/454 GS-FLX, computer analysis programs are required to assist in initial sequence data management. The data analysis softwares provided by the manufacturer of Roche/454 GS-FLX system, such as GS Reference Mapper Software or GS Amplicon Variant Analyzer Software, are not directly suitable for our purposes. Therefore, in this study, four sets of ASP16 analysis computer programs were specifically created and configured for basic analysis of HPV16 sequence reads from ASP16 experiments. The programs can deliver sorted and edited sequence alignments in FASTA format, according to their sample IDs and primer groups. The output alignments provide a direct and simple mean to investigate HPV16 sequences and to detect HPV16 integration junctions. Furthermore, the programs also deliver basic statistical data such as sequence read counts for each sorted category. The programs had been used successfully to analyze the sequence data in ASP16 experiments.

Four ASP16 experiments had been performed until now. The first two experiments, ASP16-1 and ASP16-2, were conducted by Bo Xu (Xu, 2010). It had been demonstrated that the ASP16 strategy could successfully identify HPV16 integration junctions in cell lines and clinical samples. The average sequence read lengths of these two experiments were 116 nt and 92 nt, respectively. In ASP16-2, 16 combinations of HPV16 primers were used, each located about 200 bp apart, resulting in only about 50% sequence coverage of the analyzed HPV16 E1-E2 region. This reduced the chance of finding integration junction by 50%. Therefore, to obtain better performance, the ASP16 strategy was further optimized during the two experiments ASP16-3 and ASP16-4, performed in this study. The HPV16 primer combinations were optimized to locate about 100 bp apart by increasing the number of primer combinations (Figure 2.33). The average sequence read lengths were increased from 92 nt in ASP16-2 to 105 nt and 108 nt in ASP16-3 and ASP16-4 (Table 2.15). Closer inspection into the sequence data revealed that more longer reads were obtained in ASP16-3 and ASP16-4 (Figure 2.36 and Figure 2.37). The reduction of distance between primers to 100 nt and the increase of average read length to 105-108 nt, together, contributed to 89% average sequence coverage in both experiments.

The previously determined HPV16 integration junctions of cell lines MRI-H186, MRI-H196 and SiHa were identified in the sequence reads of ASP16-3 and ASP16-4 (Table 2.18). In ASP16-2, it was not possible to identify the junctions in SiHa due to unsuitable primer locations, and the identified junction in MRI-H196 was too short to be mapped to

chromosome 11 by Blast (Xu, 2010). The verified HPV16 integration junctions in two cervical carcinoma samples CA-07C381 and CA-07C368, previously analyzed in ASP16-2, were redetected in ASP16-4. Altogether, these results demonstrate higher efficiency of the optimized ASP16 strategy in identifying HPV16 integration junctions. Cell line CaSki was also analyzed in ASP16-4. It is known to contain ~600 integrated HPV16 copies in tandem repeats, but only one 3' junction. No HPV16 integration junction could be detected in any sequence reads of CaSki in ASP16-4. This indicates that the ASP16 strategy does not allow detection of an HPV16 integration junction in a background of about 600 full-length HPV16 genomes. This might become possible if the number of sequence reads per primer is increased by two or three folds.

In ASP16-3, false-positive HPV16 integration junctions were identified in 1 and 2 sequence reads of samples CIN2/3-4242 and CIN2/3-1503, respectively. False-positive integration junctions were also identified in 12 clinical samples analyzed in ASP16-1 and ASP16-2 (Xu, 2010). These false-positive sequence reads are artifacts that may be produced during the multiple amplification steps for preparation of HPV16 amplicons. However, the number of false viral-cellular sequence reads in ASP16-3 and ASP16-4 is very low, and has decreased from the first two ASP16 experiments, possibly by optimization of the ASP16 conditions. This improvement indicates that the optimized ASP16 strategy can deliver highly specific sequence reads.

Among 25 DNA samples analyzed in ASP16-3 and ASP16-4, genuine HPV16 integration junctions have been identified in 3 out of 4 cell lines (75%) and in 6 out of 21 clinical samples (29%) (Table 2.33). All clinical samples were expected to contain integrated HPV16 DNA based on their percent integration values determined by E2/E6 qPCR. However, integration junctions could only be identified in 29% of these samples. This low determination rate might be due to a low efficiency of the ASP16 strategy, but other reasons are more likely. Closer inspection of the 21 clinical samples showed that the integration junctions were identified in 6 out of 14 samples of European lineage (43%), 0 out of 3 samples of NA1 lineage, and 0 out of 4 samples of Af1/Af2 lineages. It had been demonstrated that nucleotide variations in non-European lineages at the binding regions for E2/E6 qPCR primers and probes could result in faulty estimation of HPV16 genome status (Jiang et al, 2009). Thus, the percent integration values determined by E2/E6 qPCR of the 7 clinical samples of NA1 and Af1/Af2 lineages may be too high misrepresenting

HPV16 integration. Several studies have been performed in the past aimed at determining the proportion of CIN2/3 lesions and carcinomas harboring integrated HPV16 DNA (Briolat et al, 2007; Hudelist et al, 2004; Li et al, 2008; Saunier et al, 2008; Woodman et al, 2007). The results were different, ranging from about 20% to almost 80%, giving an average value of about 50%. The proportion of samples with integrated HPV16 DNA in our collection, 29% for all and 43% for the European lineage, fits perfectly into this range. Therefore, it seems possible that the samples, for which ASP16 failed to identify any integration junction, indeed harbor episomal HPV16 DNA.

Disruption or deletion of E2 caused by HPV16 DNA integration is thought to be a central step in cervical cancer development. Nevertheless, it is evident that some cervical carcinomas definitely contain only episomal HPV16 DNA (Hudelist et al, 2004; Park et al, 1997; Vinokurova et al, 2008). It is not yet clear whether two pathways of HPV16-induced carcinogenesis exist, one with HPV16 DNA integration as an essential step, and the other for which the integration is not obligatory because E2 inactivation is achieved by other means. It was suggested that the E2 variations of AA lineage where E2 is retained may be an alternative mechanism for increasing E6/E7 oncogene expression (Casas et al, 1999). This is supported by another study in which E2 of AA lineage was found to be impaired for repression of E6/E7 oncogene transcription (Ordóñez et al, 2004). Because the nucleotide variations in ORF E2 of NA1, Af1 and Af2 lineages are highly similar to that of AA lineage (Table 2.29 and 2.31), it is likely that the E2 of these three lineages may possess similar properties as the AA-E2, regarding the regulation of E6/E7 expression.

Overall, these results demonstrate that the optimized ASP16 strategy has the potential to analyze series of HPV16-positive clinical samples in parallel, detect HPV16 integration sites, and provide high coverage of HPV16 E1-E2 nucleotide sequences. The sequences of E1-E2 region can be used to identify mutations and assign HPV16 variants as demonstrated in this study. The integration junction analysis of the four cell lines shows that the ASP16 has a great potential to identify existent HPV16 integration junctions with the exception of samples containing a high background of full-length HPV16. The ASP16 strategy is the first method that combines the next generation sequencing technologies with HPV sequence analysis. In future development of the ASP16 strategy, it is planned to apply this strategy to other hr-HPV types, in particular HPV18 which is detected as

integrated DNA in almost 100% cervical carcinomas (Hudelist et al, 2004; Vinokurova et al, 2008). Furthermore, the Roche/454 GS-FLX will be replaced by Illumina/Solexa sequencing system, which probably has the advantage in delivering higher output volumes (~600 folds) with comparable costs. The increased output volumes will allow more samples to be analyzed per sequencing round, e.g. increasing from 20 samples to 40-50 samples, and possibly will allow the detection of integration junction in samples with high full-length HPV backgrounds, such as CaSki.

3.5 Possible consequences of HPV16 integration in samples HSIL-75857 and CIN2/3-1801

The integration of hr-HPV DNA into cellular DNA results in stable expression of viral oncogenes E6/E7 (Jeon & Lambert, 1995), which could result in transformation of the infected cells and progression into cancer (Munger et al, 1989). Furthermore, the integrated hr-HPV DNA may activate cellular proto-oncogenes or inactivate cellular tumor suppressor genes via insertional mutagenesis as suggested and demonstrated in earlier studies (Couturier et al, 1991; Peter et al, 2006; Reuter et al, 1998; Xu, 2010). When HPV DNA is integrated upstream or within and in the same orientation as a cellular gene, viral-cellular fusion transcripts can be produced which are initiated at the early viral promoter and contain the viral E6/E7 oncogenes followed by complete or truncated cellular genes (Kraus et al, 2008; Reuter et al, 1998; Wentzensen et al, 2004; Xu, 2010; Ziegert et al, 2003). The flanking cellular DNA may have an enhancer effect on viral oncogene expression, and vice versa, the HPV enhancer elements can lead to activation of cellular genes in the integration locus (von Knebel Doeberitz et al, 1991).

In sample HSIL-75857, HPV16 was found integrated into chromosome 6 in an intron of the AHI1 gene in opposite direction, and in proximity to two other cancer-relevant cellular genes, the MYB proto-oncogene and the PDE7B gene (Figure 2.53). AHI1 was suggested to be involved in tumor formation by interaction with other oncogenes such as c-myc or the tumor suppressor gene NF1 (Jiang et al, 2002). Mutations of AHI1 were reported to cause Joubert syndrome (JS) related disorders (Valente et al, 2006). The integrated HPV16 DNA may activate AHI1 gene expression. The presence of the integrated HPV16 DNA may also interfere with correct splicing of AHI1 transcripts.

MYB is a cellular proto-oncogene that can be activated by several mechanisms, including overexpression (Ramsay & Gonda, 2008). MYB was found to be overexpressed in cervical cancers and can transactivate HPV16 E6/E7 oncoprotein expression (Nurnberg et al, 1995). Integration of HPV16 DNA in the vicinity of the c-myc proto-oncogene has been shown to increase myc expression (Couturier et al, 1991; Peter et al, 2006; Xu, 2010). Thus, it is possible that the integrated HPV16 DNA, located about 181 kb downstream of myb gene, could activate the myb proto-oncogene expression, and thereby, contributing to carcinogenesis. The integrated HPV16 DNA is also located ~451 kb upstream of the cellular PDE7B gene. PDE7B was reported to be upregulated in chronic lymphocytic leukemia (CLL) (Zhang et al, 2008). Despite the relatively long distance, the integrated HPV16 DNA may affect also PDE7B expression.

In sample CIN2/3-1801, HPV16 was found integrated into chromosome 19 into an intron of the cellular REEP6 gene in opposite direction, and in proximity of two other cancer-relevant cellular genes, the APC2 and PCSK4 genes (Figure 2.54). REEP6 was proposed to be a tumor suppressor gene, and its polymorphisms are associated with colon cancer and inflammatory bowel disease (IBD) (Wellmann et al). As in the case of AHI1 gene described above, the integrated HPV16 may interfere with REEP6 expression. In addition, flanking cellular DNA may affect viral E6/E7 oncogene expression. PCSK4 is a multifunctional protein, which can activate p53 function (Batta & Kundu, 2007) as well as nonhomologous end joining and double strand break (DSB) repair activity (Batta et al, 2009). Since the integrated HPV16 DNA is located just 3998 bp upstream of PCSK4 gene in the same orientation, it is possible that viral-cellular fusion transcripts may be produced. APC2 is homologous to the tumor suppressor gene APC (van Es et al, 1999), which is associated with colon cancer (Munemitsu et al, 1995), and therefore may share functions with APC. APC2 was also suggested to be a potential tumor suppressor in ovarian cancer (Jarrett et al, 2001). The integrated HPV16 DNA may act as cis-acting regulatory sequence.

The HPV16 integration sites of both samples HSIL-75856 and CIN2/3-1801 locate in the vicinity of cellular proto-oncogenes or tumor suppressor genes. The hypothesized consequences of HPV16 integration upon the cellular genes can be investigated on RNA and protein levels, if the respective materials are available from the clinical samples. The HPV16 integration sites of the two samples support the assumption that HPV integration

is important for cervical carcinomas not only by deregulating/activating E6/E7 oncogene expression, but also by altering cancer-relevant cellular genes.

Up to the present, more than 200 HPV integration loci had been identified (Kraus et al, 2008; Wentzensen et al, 2004; Yu et al, 2005). However, the consequences of viral integration have not been investigated for these identified integration loci. With the vast resources of human genome sequences, genes and proteins now available through many database providers, such as NCBI, it is possible to create computer algorithms that can predict the outcomes of HPV integrations. Many tools are already available, such as gene structure and alternative splicing prediction (Coward et al, 2002), regulatory element prediction (Robertson et al, 2006) and protein interaction prediction (Jansen et al, 2003). In combining several available tools, different aspects of the effects of HPV integration could be virtually investigated.

4. Materials and methods

4.1 Chemicals, commercial media/solutions and antibiotics

Acetic acid, glacial	Sigma Aldrich, Munich
Agarose	Sigma Aldrich, Munich
Ampicillin, sodium salt, powder	AppliChem, Darmstadt
Bacto agar	Becton Dickinson, Heidelberg
Bovine serum albumin (BSA)	New England Biolabs, USA
Bromophenol blue, sodium salt	Serva, Heidelberg
Chloroform	Roth, Karlsruhe
deoxyadenosine triphosphate (dATP)	Roche Applied Science, Mannheim
deoxycytidine triphosphate (dCTP)	Roche Applied Science, Mannheim
deoxyguanosine triphosphate (dGTP)	Roche Applied Science, Mannheim
deoxythymidine triphosphate (dTTP)	Roche Applied Science, Mannheim
Diethylpyrocarbonate (DEPC)	AppliChem, Darmstadt
Dulbecco's Modified Eagle's Medium (DMEM)	Sigma Aldrich, Munich
Dimethyl sulfoxide (DMSO)	AppliChem, Darmstadt
Ethylenediaminetetraacetic acid (EDTA), disodium salt dihydrate	Sigma Aldrich, Munich
Ethanol, absolute	Sigma Aldrich, Munich
Fetal bovine serum (FBS), superior	Biochrom, Berlin
Ficoll 400	Serva, Heidelberg
Formaldehyde	Merck, Darmstadt
Formamide	Merck, Darmstadt
Glucose	Sigma Aldrich, Munich
Glycerol	Sigma Aldrich, Munich
Hydrochloric acid (HCl), 37% (w/v)	Sigma Aldrich, Munich
Isoamyl alcohol	Sigma Aldrich, Munich
Isopropanol	Sigma Aldrich, Munich
Kanamycin sulfate, powder	AppliChem, Darmstadt
Magnesium chloride (MgCl ₂)	Sigma Aldrich, Munich
Magnesium sulfate (MgSO ₄)	Sigma Aldrich, Munich
NEN [α - ³² P] dCTP	PerkinElmer, Rodgau-Jügesheim
Penicillin/streptomycin solution	Invitrogen, Karlsruhe
Phenol	Roth, Karlsruhe
Polyvinylpyrrolidone	Sigma Aldrich, Munich

Potassium chloride (KCl), powder	Sigma Aldrich, Munich
Potassium phosphate monobasic (KH ₂ PO ₄)	Sigma Aldrich, Munich
Random hexamer - 5'NNNN*N*N3' (phosphothioate modified)	Thermo Fisher Scientific, Ulm
Redivue [α - ³² P] dCTP	GE Healthcare, Munich
Sarkosyl, sodium salt	Sigma Aldrich, Munich
Sodium acetate, powder	Sigma Aldrich, Munich
Sodium chloride (NaCl), powder	Sigma Aldrich, Munich
Sodium citrate dihydrate	Sigma Aldrich, Munich
Sodium dodecyl sulfate (SDS)	Sigma Aldrich, Munich
Sodium hydroxide (NaOH)	Sigma Aldrich, Munich
Sodium phosphate dibasic (Na ₂ HPO ₄)	Sigma Aldrich, Munich
t-RNA	Sigma Aldrich, Munich
Trizma base, powder	Sigma Aldrich, Munich
Trypsin/EDTA solution	Invitrogen, Karlsruhe
Tumor necrosis factor alpha (TNFalpha), recombinant	Strathmann Biotec, Hamburg
UltraPure water, DNase/RNase-free	Invitrogen, Karlsruhe
Tryptone	Becton Dickinson, Heidelberg
Yeast extract	Becton Dickinson, Heidelberg
Xylene cyanol FF, sodium salt	Serva, Heidelberg

4.2 Buffers, stock solutions and media

6xBPB DNA loading buffer

0.25% (w/v)	bromophenol blue
30% (v/v)	glycerol

6xXC DNA loading buffer

0.25% (w/v)	xylene cyanol
30% (v/v)	glycerol

50xTAE electrophoresis buffer (for 1 L)

242 g	Trizma
57.1 ml	glacial acetic acid
100 ml	0.5 M EDTA, pH 8.0

10xTE stock solution (for 1 L)

100 ml	1 M Tris-HCl, pH 8.0
20 ml	0.5 M EDTA, pH 8.0
880 ml	deionized water

1xTE/RNase buffer (for 100 ml)

10 ml	10xTE
1 ml	10 mg/ml RNase A
89 ml	deionized water

TEG buffer

50 mM	glucose (filtered)
10 mM	EDTA, pH 8.0
25 mM	Tris-HCl

NaOH/SDS solution

0.2 M	NaOH
1%	SDS

20xSSC stock solution, pH 7.0

3 M	NaCl
0.3 M	Sodium citrate

10xTNE stock solution

100 mM	Trizma base, pH 8.0
1 M	NaCl
100 mM	EDTA, pH 8.0
2%	SDS

50xDenhardt's solution

1% (w/v)	Ficoll 400
1% (w/v)	polyvinylpyrrolidone
1% (w/v)	BSA

Hybridization solution

50 mM	sodium phosphaste, pH 6.4
1%	SDS
1x	Denhardt's solution
5x	SSC
0.1 mg/ml	t-RNA

LB medium with selective antibiotics (for 1 L)

10 g	NaCl
10 g	tryptone
5 g	yeast extract

Adjusted to pH 7.0. After autoclaved, one of the following antibiotics was added:

1 ml	100 mg/ml ampicillin
or 5 ml	10 mg/ml kanamycin

LB agar (1 L) with selective antibiotics

1 L LB medium without antibiotics

14 g bacto agar

After autoclaved, one of the following antibiotics was added:

1 ml 100 mg/ml ampicillin

or 5 ml 10 mg/ml kanamycin

SOC medium (1 L)

0.5 g NaCl

20 g tryptone

5 g yeast extract

20 ml 1 M glucose

Filter-sterilized.

10xPBS stock solution (1 L), pH 7.4

80 g NaCl

2 g KCl

14.4 g Na₂HPO₄

2.4 g KH₂PO₄

PBS/EDTA solution

1x PBS

2 M EDTA

3x lysis buffer

2x TE

3% (w/v) sarkosyl

Complete DMEM medium (for cell culture) with antibiotics

500 ml DMEM

50 ml FBS

5 ml penicillin/streptomycin (50000 units each)

Cell freezing medium (for 10 ml)

8 ml complete DMEM medium

1 ml FBS

1 ml DMSO

CIA (Chloroform:Isoamyl alcohol = 24:1)

24 part (v/v) Chloroform

1 part (v/v) Isoamyl alcohol

Phenol* (Phenol:CIA = 1:1)

1 part (v/v) CIA

1 part (v/v) Phenol

5xAnnealing buffer (for phi29 DNA polymerase)

85 mM	Tris-HCl, pH 8.0
85 mM	MgCl ₂

4.3 Commercial kits and enzymes

Calf Intestine Alkaline Phosphatase (CIAP)	Fermentas Life Sciences, Canada
Expand long template PCR system	Roche Applied Science, Mannheim
GenomePlex WGA2 kit	Sigma Aldrich, Munich
Pfu DNA Polymerase	Promega, USA
Phi29 DNA Polymerase	Fermentas Life Sciences, Canada
Proteinase K	Merck, Darmstadt
Qiagen Multiplex PCR kit	Qiagen, Hilden
Qiagen Plasmid Maxi kit	Qiagen, Hilden
QIAquick Gel Extraction kit	Qiagen, Hilden
QIAquick PCR Purification kit	Qiagen, Hilden
QIAprep Spin Miniprep kit	Qiagen, Hilden
Rediprime II DNA Labeling System	GE Healthcare, Munich
RNase A	Qiagen, Hilden
StrataClone PCR Cloning kit	Stratagene, La Jolla, USA
StrataClone Blunt PCR Cloning kit	Stratagene, La Jolla, USA
T4 DNA Ligase (2000 U/ml)	New England Biolabs, USA
Taq DNA Polymerase kit	Invitrogen, Karlsruhe
TOPO TA PCR Cloning kit	Invitrogen, Karlsruhe
TOPO Blunt PCR Cloning kit	Invitrogen, Karlsruhe

4.4 Laboratory equipments and commercial materials

BioMag streptavidin beads	Polysciences, Heidelberg
Biomax MR film	GE Healthcare, Munich
E-Gel iBase Power System	Invitrogen, Karlsruhe
E-Gel Safe Imager Transilluminator	Invitrogen, Karlsruhe
E-Gel SizeSelect	Invitrogen, Karlsruhe
Hybond-N+ nylon membrane	GE Healthcare, Munich
MicroSpin G-25 column	GE Healthcare, Munich
Sterile filter (0.45 µm)	Millipore, Schwalbach/Ts.
Thermocycler PTC-200	MJ Research, Waltham, USA
Whatman 3MM paper	Schleicher & Schuell, Dassel

4.5 Computer software and Internet resources

BLAST version 2.2.17

<ftp://ftp.ncbi.nlm.nih.gov/blast/executables/release/>

BLAST executables are based on NCBI C Toolkit and consist of several stand-alone executable programs, including *blastall* and *formatdb*. In this study, *formatdb* was used to create sequence databases on the analyzing computer, and *blastall* was used for blasting sequence reads of ASP16 (as queries) to the sequences in the created databases (as subjects).

ClustalW version 2.0.10

<http://www.clustal.org/>

ClustalW is the command-line version of Clustal, a multiple sequence alignment program. It accepts and delivers several input/output sequence formats. FASTA format was used in this study for both inputs and outputs.

Geneious Pro version 4.8.5

<http://www.geneious.com/>

Geneious Pro is an integrated, cross-platform bioinformatics software suited for manipulating, finding, sharing, and exploring biological data. It was used mainly for sequence visualization and edition.

HUSAR: Heidelberg Unix Sequence Analysis Resources

<http://genome.dkfz-heidelberg.de/>

HUSAR is a program package providing access to more than 200 databases and offering more than 260 applications for DNA and protein analysis, including many widely known programs such as *clustal* (for sequence alignment), *bl2Seq* (for sequence comparison) and *primer* (for primer design).

NCBI: National Center for Biotechnology Information

<http://www.ncbi.nlm.nih.gov/>

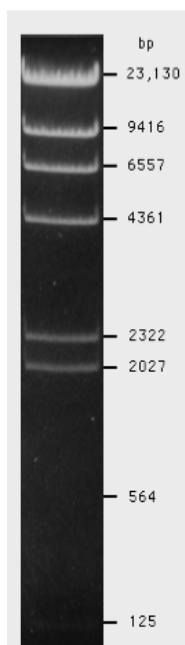
NCBI contains various databases of biotechnology information, including nucleotide and protein sequence databases of different organisms. It was used for searching, comparing and retrieving sequences.

For nucleotide sequences of the twenty-four human chromosomes, the reference sequences with accession numbers NC_000001, NC_000002,..., NC_000024 are used in this work. Each accession number comes with a version number, identified with suffix “.(version-number)”.

Perl application version 5.8.8

<http://www.perl.org/>

Perl version 5.8.8 is a program that reads a Perl language program, translates it into instructions the computer can understand, and runs (or executes) it. The developed ASP16 data analysis programs (see Results section 2.2.1) were executed in Perl 5.8.8 under Mac operating system.

4.6 DNA molecular markers**Lambda DNA/HindIII marker**

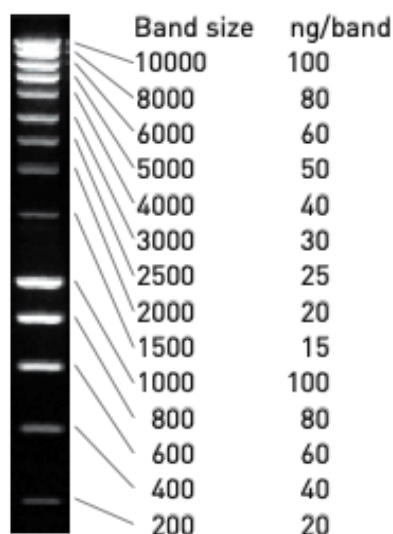
This DNA ladder is frequently used when running gDNA because of its high molecular weight bands

Practical size range: 2 – 23 kb

0.5 µg/lane

1.2% agarose gel, stained with ethidium bromide

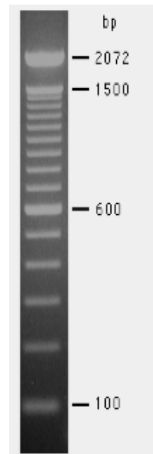
(Invitrogen, Karlsruhe)

SmartLadder

Size range: 200 bp – 10 kb

5 µl/lane (720 ng)

(Eurogentec, Köln)

100bp Ladder

Size range: 100 bp – 2 kb

0.5 µg/lane

2% agarose gel, stained with ethidium bromide

(Invitrogen, Karlsruhe)

4.7 Oligonucleotide primers

Oligonucleotide primers for routine PCRs were designed so that they have optimal melting temperature about 60°C, and are 18-24 nt in length.

4.7.1 Primers for cellular genes

The primers, shown in Table 4.1, were used to amplify parts of cellular genes after DNA preparation to determine the integrity of the DNA template. β -globin and GAPDH are house-keeping genes, while c-myc is a proto-oncogene.

Table 4.1: Oligonucleotide primers for cellular genes.

Primer name*	Sequence (5' to 3')	Target cellular gene	Expected product size
β -globin-F	ACACAACGTGTTCCTACTAGC	β -globin	110 bp
β -globin-R	CAACTTCATCCACGTTCCACC		
GAPDH-F	CACCACCAACTGCTTAGCAC	GAPDH	367 bp
GAPDH-R	GAGGCAGGGATGATGTTCTG		
c-myc-F	CTTTATAATGCGAGGGTCTGG	c-Myc	591 bp
c-myc-R	GCTTACCTGGTTTCCACTACC		

* Binding directions are indicated. F: forward. R: reverse.

4.7.2 Primers for HPV68 analysis

The primers used for amplification of HPV68 DNA described in Results section 2.1 are shown in Table 4.2. The primers with names starting with “ME” were used mainly for amplification of the integrated HPV68b in cell lines ME180 and ME180R. The primers with names starting with “H68” were used mainly for analysis of HPV68b DNA in cervical samples.

Table 4.2: Oligonucleotide primers for HPV68.

Primer name	Sequence (5' to 3')	Primer binding sites *
ME-01 §	GGAATTATGCCTATTTGACA	pos. 45577909-F (NC_000018.9), upstream of the integrated HPV68b in ME180/ME180R
ME-04	TGGTATTTTGGTGTGGTTTGTG	pos. 3860-F of HPV68b-CIN2 genome
ME-874 §	ACACACCAATGGAGAGGAGTG	pos. 45578156-F of NC_000018.9, upstream of the integrated HPV68b in ME180/ME180R
ME-1349	TAGTCACAGGTGCAACCAC	pos. 7208-R of HPV68b-CIN2 genome
ME-6651	TTTGTGTGTCCGTGGTGTG	pos. 782-F of HPV68b-CIN2 genome
ME-6805	CGCGACAGATACAGGTTTCAG	pos. 940-F of HPV68b-CIN2 genome
ME-7094	CATTTTCCCTTCATCTCC	pos. 1289-R of HPV68b-CIN2 genome
ME-7718	ACCATTGGAGGTCTTTGCTG	pos. 3984-F of HPV68b-CIN2 genome
ME-7810	ACATGCACACACAAAACCAC	pos. 3892-R of HPV68b-CIN2 genome
ME-7972 §	TCTGGCTCAGGAAATCGC	pos. 45578405-R of NC_000018.9, downstream of the integrated HPV68b in ME180/ME180R
H68-128R	TCTGTCCGTGTAGTTGCCTTC	Primer positions and binding directions on CIN2-HPV68b** genome are indicated in the primer names, after the prefix "H68-".
H68-377R	CTTCGTTTTGAATTTAGGTGCC	
H68-782F	TTTGTGTGTCCGTGGTGTG	
H68-830R	TGGCCATTGCAGATTACTGG	
H68-1107F	GCAGCCCTTTAGCAAAGTC	
H68-1175R	GTTGTCTTGCTGTGTACTGC	
H68-1709F	CTTTTGCAGCCACCAAAATT	
H68-1728R	AATTTTGGTGGCTGCAAAAG	
H68-2915F	GATGGCACTAGAGAGCATTCG	
H68-2920F	GCCATCTGCAGTTCAATGG	
H68-4145F	AAGCGTGCATCTGCAACTG	
H68-4572R	ACTTGCACAGACCCAGACGAAG	
H68-5629F	GGTTATTAAGTGTAGGCCATCC	
H68-5877F	AATAGGCTAGATGATACTGAG	
H68-5931R	CCTTAGGATTTTGTGGAGGAAAC	
H68-6479R	ACAAATACCATTTGTTGTGTCCC	
H68-7192F	GTTGCACCCTGTGACTAACATATG	
H68-7954F	TGCCTAATAGCATAGTTGGCC	
MY11	GCMCAGGGWCATAAATGG M=A/C, W=A/T, Y=C/T, R=A/G	pos. 6453-F of HPV68b-CIN2 genome

* Binding positions (5'end) and binding directions are indicated. F: forward. R: reverse.

** See Results section 2.1.4. The sequence of HPV68b-CIN2 is shown in Appendix A2.

§ Primers bind to human chromosome 18, flanking the integrated HPV68b in ME180 and ME180R.

4.7.3 Primers for HPV16 and HPV16 integration junctions

The primers used for amplification of HPV16 DNA and HPV16 integration junctions described in Results section 2.2 are listed in Table 4.3.

Table 4.3: Oligonucleotide primers for HPV16 and HPV16 integration junctions.

A: Primers specific for HPV16 DNA.

Primer name *	Sequence (5' to 3')	Remarks
H16-44F	GGTTGAACCGAAACCGGTTAG	
H16-691R	GTCCAGCTGGACCATCTATTTTC	
H16-1064F	CACATGCGTTGTTTACTGCAC	binds at the same pos. as primer E09 (1064) **
H16-1105F	GATGCAGTACAGGTTCTAAAACG	binds at the same pos. as primer BE-19 **
H16-1115F	AGGTTCTAAAACGAAAGTATTGG	binds at the same pos. as primer E19 (1115) **
H16-1656F	TCAAAGTTTAGCATGTTTCATGGG	
H16-2174F	GTGATTGGAAGCAAATTGTTATG	binds at the same pos. as primer E12 (2174) **
H16-2403F	AGCAGATGCCAAAATAGGTATG	binds at the same pos. as primer E05 (2403) **
H16-3121F	TGGAGACATATGCAATACAATGC	binds at the same pos. as primer E07 (3121) **
H16-3339F	GTCTACATCTGTGTTTAGCAGC	

* Binding positions (5'end) and binding directions are indicated after prefix "H16". F: forward. R: reverse.

** Primers are shown in Table 4.4 and Table 4.5.

(Table 4.3, continued)

B: Primers specific for human chromosomes.

Primer name	Sequence (5' to 3')	Primer binding sites ***
NC06-R-1.1	AGAAAACCTCCCTAGTCTGAAATC	chrom.6 pos. 097896916-F (NC_000006.11)
NC06-R-1.2	TCAACTAGGATTACCTACCCTCC	chrom.6 pos. 097896741-F (NC_000006.11)
NC06-R-2.1	GTGTTTTTCAATTATTTGGAAAGG	chrom.6 pos. 135721349-R (NC_000006.11)
NC06-R-2.2	ATTGTAGTCGAATCCACCTCGG	chrom.6 pos. 135721494-R (NC_000006.11)
NC07-R-1.1	AGTGGTGTCCCCGCTGTAAG	chrom.7 pos. 72454604-R, 72818829-R and 75010302-F (NC_000007.13)
NC11-R-1.1	GAGGAGTGACTTTGAATAGAATGG	chrom.11 pos. 123658930-F (NC_000011.9)
NC11-R-1.2	AGTGGGGAGGTAGCTTCTGG	chrom.11 pos. 123658817-F (NC_000011.9)
NC19-R-1.1	GGGGCTGATGGCATGGAGTG	chrom.19 pos. 01494382-F (NC_000019.9)
NC19-R-1.2	GTTTGCCTCCAAGACAG	chrom.19 pos. 01494273-F (NC_000019.9)

*** Binding positions (5'end) and binding directions are indicated. F: forward. R: reverse.

4.7.4 HPV16 primers for ASP16 strategy

In the ASP16 strategy, two sets of HPV16 forward primers and one set of reverse primers are used in the HPV16 DNA enrichment step (see section 4.22.2). The first set of HPV16 forward primers are shown in Table 4.4. These primers are labeled with biotin at the 5'end, and were used in linear amplification step. The second set of HPV16 forward primers are shown in Table 4.5. They are bipartite nested HPV16 primers that contain the 19-nt Roche-A sequence (RA) at the 5'end. The reverse primers, shown in Table 4.6, are tripartite. They contain the 19-nt Roche-B sequence (RB), followed by a 4-nt barcode and the 18-nt GenomePlex universal adapter sequence (GPUA). These reverse primers were used in combination with the bipartite nested HPV16 primers for the nested HPV16 PCR, following the linear amplification step in ASP16 strategy.

Table 4.4: Biotin-labelled HPV16 forward primers for linear amplification in ASP16 strategy.

Short name	Full name*	Sequence (5' to 3')
BE-01	5B-810F	AATGGGCACACTAGGAATTGTG
BE-02	5B-1261F	GGGTATGGCAACTACTGAAGTGG
BE-03	5B-1562F	GTTGCGATTGGTGTATTGCTG
BE-04	5B-1938F	ACAGATGGTACAATGGGCCTAC
BE-05	5B-2389F	TGGTTACAACCAATTAGCAGATGC
BE-06	5B-2705F	CAAGGACGTGGTCCAGATTAAG
BE-07	5B-3101F	ACAGTGGAAGTGCAGTTTGATG
BE-08	5B-3542F	GACAGTGCTCCAATCCTCACTG
BE-09	5B-1046F	AGGCAGAAACAGAGACAGCAC
BE-10	5B-1394F	GAGAGGGTGTTAGTGAAAGACAC
BE-11	5B-1760F	GTGTGTCTCCAATGTGTATGATG
BE-12	5B-2151F	ATGTGATAGGGTAGATGATGGAG
BE-13	5B-2539F	TGCCCTCCATTATTAATTACATC
BE-14	5B-2912F	CATATTAACCACCAAGTGGTGC
BE-15	5B-3318F	CGGGTGGTGAGGTAATATTATG
BE-16	5B-3762F	CATATGATAGTGAATGGCAACG
BE-17	5B-790F	CGTACTTTGGAAGACCTGTTAATG
BE-18	5B-952F	GGGGATGCTATATCAGATGACG
BE-19	5B-1105F	GATGCAGTACAGTTCTAAAACG
BE-20	5B-1336F	AGTCAGTATAGTGGTGGAGTGG
BE-21	5B-1457F	TAAAACTAGTAATGCAAAGGCAG
BE-22	5B-1653F	CATTCAAAGTTTAGCATGTTTCATG
BE-23	5B-1849F	ATTAGTGAAGTGTATGGAGACACG

(Table 4.4, continued)

Short name	Full name*	Sequence (5' to 3')
BE-24	5B-2053F	TCACAGGCAGAAAAATTGTAAAGG
BE-25	5B-2277F	TTGCATATTACTATATGGTGACG
BE-26	5B-2461F	GATGACAATTTAAGAAATGCATTG
BE-27	5B-2613F	GGTGGTGTTTACATTTCCTAATG
BE-28	5B-2842F	TATAGACTATTGGAAACACATGCG
BE-29	5B-3021F	ATGAAAAGTGGACATTACAAGACG
BE-30	5B-3189F	TAACTGTGGTAGAGGGTCAAGTTG
BE-31	5B-3444F	GCACCGAAGAAACACAGACG
BE-32	5B-3680F	TGTACATTGTATACTGCAGTGTCTG

* Binding positions (5' end) and binding direction are indicated after prefix "5B". F: forward.

Table 4.5: Nested HPV16 forward primers for semi-nested PCR in ASP16 strategy.

Short name	Full name*	Sequence (5' to 3')
E01	RA-857F	GCCTCCCTCGCGCCATCAG AATCTACCATGGCTGATCCTG
E02	RA-1275F	GCCTCCCTCGCGCCATCAG TGAAGTGGAACTCAGCAGATG
E03	RA-1576F	GCCTCCCTCGCGCCATCAG ATTGCTGCATTGGACTTACAC
E04	RA-1951F	GCCTCCCTCGCGCCATCAG TGGGCTACGATAATGACATAG
E05	RA-2403F	GCCTCCCTCGCGCCATCAG AGCAGATGCCAAAATAGGTATG
E06	RA-2723F	GCCTCCCTCGCGCCATCAG TAAGTTTGACAGAGGACGAG
E07	RA-3121F	GCCTCCCTCGCGCCATCAG TGGAGACATATGCAATACAATGC
E08	RA-3555F	GCCTCCCTCGCGCCATCAG TCCTCACTGCATTAAACAGCTC
E09	RA-1064F	GCCTCCCTCGCGCCATCAG CACATGCGTTGTTTACTGCAC
E10	RA-1411F	GCCTCCCTCGCGCCATCAG AGACACACTATATGCCAAACACC
E11	RA-1785F	GCCTCCCTCGCGCCATCAG AGAGCCTCCAAAATTGCGTAG
E12	RA-2174F	GCCTCCCTCGCGCCATCAG GTGATTGGAAGCAAATTGTTATG
E13	RA-2569F	GCCTCCCTCGCGCCATCAG AATGCTGGTACAGATTCTAGGTG
E14	RA-2933F	GCCTCCCTCGCGCCATCAG CCAACACTGGCTGTATCAAAG
E15	RA-3339F	GCCTCCCTCGCGCCATCAG GTCCTACATCTGTGTTTAGCAGC
E16	RA-3778F	GCCTCCCTCGCGCCATCAG GCAACGTGACCAATTTTTGTC
E17	RA-810F	GCCTCCCTCGCGCCATCAG AATGGGCACACTAGGAATTGTG
E18	RA-968F	GCCTCCCTCGCGCCATCAG ATGACGAGAACGAAAATGACAG
E19	RA-1115F	GCCTCCCTCGCGCCATCAG AGGTTCTAAAACGAAAGTATTTGG
E20	RA-1361F	GCCTCCCTCGCGCCATCAG GTGGTTGCAGTCAGTACAGTAGTG
E21	RA-1471F	GCCTCCCTCGCGCCATCAG GCAAAGGCAGCAATGTTAGC
E22	RA-1672F	GCCTCCCTCGCGCCATCAG TCATGGGGAATGGTTGTGTTAC
E23	RA-1860F	GCCTCCCTCGCGCCATCAG GTATGGAGACACGCCAGAATG
E24	RA-2069F	GCCTCCCTCGCGCCATCAG TAAAGGATTGTGCAACAATGTG
E25	RA-2288F	GCCTCCCTCGCGCCATCAG TATATGGTGCAGCTAACACAGG
E26	RA-2476F	GCCTCCCTCGCGCCATCAG AATGCATTGGATGGAAATTTAG
E27	RA-2628F	GCCTCCCTCGCGCCATCAG TCCTAATGAGTTTCCATTGTGACG
E28	RA-2857F	GCCTCCCTCGCGCCATCAG ACACATGCGCCTAGAATGTG
E29	RA-3037F	GCCTCCCTCGCGCCATCAG ACAAGACGTTAGCCTTGAAGTG
E30	RA-3199F	GCCTCCCTCGCGCCATCAG AGAGGGTCAAGTTGACTATTATGG
E31	RA-3455F	GCCTCCCTCGCGCCATCAG ACACAGACGACTATCCAGCGAC
E32	RA-3696F	GCCTCCCTCGCGCCATCAG CAGTGTCTGTACATGGCATTG

* Binding positions (5' end) and binding direction are indicated after prefix "RA". F: forward.

Table 4.6: Reverse barcode primers for semi-nested PCR in ASP16 strategy.

Primer name	Sequence (5' to 3')
RB-B01	GCCTTGCCAGCCCGCTCAG TGAC TGTGTTGGGTGTGTTTGG
RB-B02	GCCTTGCCAGCCCGCTCAG AGAC TGTGTTGGGTGTGTTTGG
RB-B03	GCCTTGCCAGCCCGCTCAG TCAC TGTGTTGGGTGTGTTTGG
RB-B04	GCCTTGCCAGCCCGCTCAG ACAC TGTGTTGGGTGTGTTTGG
RB-B05	GCCTTGCCAGCCCGCTCAG TGTC TGTGTTGGGTGTGTTTGG
RB-B06	GCCTTGCCAGCCCGCTCAG AGTC TGTGTTGGGTGTGTTTGG
RB-B07	GCCTTGCCAGCCCGCTCAG TCTC TGTGTTGGGTGTGTTTGG

(Table 4.6, continued)

Primer name	Sequence (5' to 3')
RB-B08	GCCTTGCCAGCCCGCTCAG ACTC TGTGTTGGGTGTGTTTGG
RB-B09	GCCTTGCCAGCCCGCTCAG CTGA TGTGTTGGGTGTGTTTGG
RB-B10	GCCTTGCCAGCCCGCTCAG CAGA TGTGTTGGGTGTGTTTGG
RB-B11	GCCTTGCCAGCCCGCTCAG CTGA TGTGTTGGGTGTGTTTGG
RB-B12	GCCTTGCCAGCCCGCTCAG CACA TGTGTTGGGTGTGTTTGG
RB-B13	GCCTTGCCAGCCCGCTCAG TAGC TGTGTTGGGTGTGTTTGG
RB-B14	GCCTTGCCAGCCCGCTCAG ATGC TGTGTTGGGTGTGTTTGG
RB-B15	GCCTTGCCAGCCCGCTCAG TACA TGTGTTGGGTGTGTTTGG
RB-B16	GCCTTGCCAGCCCGCTCAG ATCA TGTGTTGGGTGTGTTTGG
RB-B17	GCCTTGCCAGCCCGCTCAG TGCA TGTGTTGGGTGTGTTTGG
RB-B18	GCCTTGCCAGCCCGCTCAG TCGA TGTGTTGGGTGTGTTTGG
RB-B19	GCCTTGCCAGCCCGCTCAG AGCA TGTGTTGGGTGTGTTTGG
RB-B20	GCCTTGCCAGCCCGCTCAG ACGA TGTGTTGGGTGTGTTTGG
RB-B21	GCCTTGCCAGCCCGCTCAG CATC TGTGTTGGGTGTGTTTGG
RB-B22	GCCTTGCCAGCCCGCTCAG CTAC TGTGTTGGGTGTGTTTGG
RB-B23	GCCTTGCCAGCCCGCTCAG CTGC TGTGTTGGGTGTGTTTGG
RB-B24	GCCTTGCCAGCCCGCTCAG CAGC TGTGTTGGGTGTGTTTGG

4.8 Cervical carcinoma cell lines

Frozen stocks of the following cervical carcinoma cell lines were provided by Elisabeth Schwarz. The integrated HPV types and the cell lines are indicated in Table 4.7.

Table 4.7: Cervical carcinoma cell lines.

Cell line name	Integrated HPV type	Source
CaSki	HPV16	American Type Culture Collection
ME180	HPV68b	American Type Culture Collection
ME180R	HPV68b	M. Schaadt. Medizinische Klinik Universität Köln, Germany
MRI-H186	HPV16	American Type Culture Collection
MRI-H196	HPV16	American Type Culture Collection
SiHa	HPV16	American Type Culture Collection

4.9 Clinical DNA samples

Our collaboration partners in Besancon and Reims provided us with the clinical DNA samples from cervical scrapes or lesions. Through the routine cervical screening programs in Besancon and Reims, cervical scrapes were collected from women of different ages using the ThinPrep system. In brief, the cervical cells were collected by a cyto-brush and suspended in the ThinPrep PreserveCyt solution where the cells are fixed. This allows long-term storage of the samples. All samples underwent cytology analysis and HPV genotyping. For selected samples, DNA was extracted using Qiagen EZ1 DNA tissue kit (Qiagen, Hilden). Selected HPV16- and HPV68-positive DNA samples were sent to us. For an example, see (Briolat et al, 2007) for a description of the sample preparation and analysis.

4.10 *In vitro* cultivation of cervical carcinoma cell lines

All procedures took place under a laminar flow hood. All culture media and solutions used were prewarmed at 37°C. Cell cultures were incubated in a cell incubator at 37°C with 5% CO₂ and 95% humidity.

The cells were taken from 1-ml cell line aliquot stocks stored in liquid nitrogen. After rapidly thawing to room temperature, the cells were transferred to 14-ml complete DMEM medium in a 75cm² T-flask, and suspended by repeat pipetting. As soon as the cells adhered to the bottom of the flask, the medium was renewed to clear the residual DMSO from the freezing medium. Medium was changed every 2-3 day. When the cells had grown to 90% confluent, they were subcultured.

For subculturing (splitting), confluent cells were rinsed with 5 ml trypsin/EDTA, then incubated with new 3 ml trypsin/EDTA at 37°C for 3-10 minutes until the adherent cells rounded up. Cell suspensions were then diluted 1:10 or 1:20 with complete DMEM medium and put in a new 75cm² T-flask. After the dilution, the cells were suspended by repeat pipetting. The cultures were cultivated in the incubator.

For long-term storage of the cells, they were cultured in 175cm² T-flasks until reaching 80% confluency. After trypsinization, the cell pellets were collected by centrifugation at 800 rpm for 10 minutes at 4°C. Cell pellets were washed with 10 ml complete DMEM medium and collected by centrifugation again. Medium was discarded. Cell pellets were resuspended in 5 ml freezing medium and transferred to five cryogenic vials. The vials were cooled down slowly to -80°C, and then stored in liquid nitrogen.

4.11 Isolation of TNFalpha-resistant cells from ME180

The cervical carcinoma cell line ME180 was cultured in complete DMEM with 5 nM and 10 nM TNFalpha (Strathmann Biotec, Hamburg) continuously for 18 weeks. During this time, the culturing media with fresh TNFalpha were changed every 2-3 days. After 18 weeks, the survival cells were pooled and cultured in complete DMEM medium without TNFalpha for another 4 weeks. The two pooled cell populations, obtained from 5 nM and 10 mM TNFalpha treatments, were designated ME180-2A and ME180-3A, respectively.

4.12 TNFalpha cytotoxicity assay

The ME180, ME180R, ME180-2A and ME180-3A cells (see Results section 2.1.3) were seeded on a 96-well culture plate at density of 5000 and 10000 cells/well; 12 wells per cell concentration per cell population. The cells were cultured in complete DMEM medium overnight, allowing them to attach to the plate. In the next day, four different concentrations of TNFalpha (0, 2.5, 5 and 10 nM) were given; 3 wells per TNFalpha concentration per cell density per cell population. The cultures were incubated for 72 hours. Afterward, the media were discarded and the cells were fixed in 4% glutaraldehyde for 3 minutes and then washed in tab water. The fixed cells were stained with 1% crystal violet for 3 minutes and then washed in tab water. The plate was left to dry completely at room temperature. Prior to the measurement, 100 µl of 33% acetic acid was added to each well. The absorbance of the crystal violet dye, directly proportional to the cell amount in each well, was measured at wavelength 560 nm. The average absorbance values of each triplicate (three wells of the same TNFalpha concentration, seeding density and cell source) were determined. The cell viability was determined, in percentage, as the absorbance value after treatments with TNFalpha (2.5, 5 and 10 nM) over the value in absence of TNFalpha (0 nM) for each cell population.

4.13 Isolation of genomic DNA from monolayer cell cultures

Genomic DNA (gDNA) was isolated from cultured cells using a standard procedure (Strauss, 2001), which involves proteinase K digestion for protein degradation, repeated extractions with phenol/chloroform/isoamylalcohol for deproteinization, ethanol precipitation, and resuspension of DNA.

Cells were cultivated until 90-95% confluent in 75cm² T-flasks. The cell monolayers were rinsed twice with cold PBS/EDTA solution. To lyse the cells, 3 ml cold PBS/EDTA and 1.5 ml 3x lysis buffer were added to the flask. The lysis buffer was applied to the monolayer, and the flask was swirled slowly until the solution became viscous. The cell lysates were transferred to a 50-ml Falcon tube and 90 µl Proteinase K (final concentration 0.2 mg/ml) was added. The lysates were incubated at 55°C overnight, allowing histones to be digested and cellular nucleases to be inactivated by the enzyme. Next day, the reactions were cooled down to room temperature before 5 ml Phenol* was added. The reaction was suspended by inverting the tubes for 30 minutes, and the aqueous and phenol phases were separated by centrifugation at 2000 rpm for 10 minutes at 20°C. The aqueous layer (4.5 ml) was transferred to a new 14-ml Falcon tube, and 4.5 ml CIA solution was added. The solutions were mixed by inverting the tube for 30 minutes, and the aqueous and phenol phases were separated by centrifugation at 2000 rpm for 10 minutes at 20°C.

The aqueous layer (4.2 ml) was transferred to a new 14-ml Falcon tube, and the gDNA was precipitated by adding 0.1 volume of 3 M sodium acetate (pH 5.2) and 2 volumes ice-cold absolute ethanol. gDNA was collected by centrifugation at 5000 rpm for 30 minutes at 4°C, then washed with 5 ml ice-cold 75% ethanol. The DNA pellet was left at room temperature for 30 minutes to clear the residue ethanol. gDNA was suspended in 2 ml 1xTE buffer and stored at 4°C.

Due to the high molecular weight of genomic DNA, their quality and concentration were determined by gel electrophoresis, in comparison with Lambda DNA/HindIII DNA marker.

4.14 Polymerase chain reaction (PCR)

Polymerase chain reaction is a molecular technique used to amplify specific piece of DNA from trace amount of target molecules (Kuslich et al, 2008). It can produce millions copies of the target DNA. In this work, three PCR systems had been used: PCR with Taq DNA polymerase, with Pfu DNA polymerase, and with the Expand Long Template PCR System kit (Roche). The PCR reactions were set up in 250- μ l PCR tubes and run in a Thermocycler PTC-200, which can handle PCR reaction volume up to 50 μ l. The specificity of the PCR products was analyzed by agarose gel electrophoresis where 1/10 of each PCR products were run together with a suitable DNA marker

PCR system 1: PCR with Taq DNA polymerase

This was used for routine PCR. There is no proof-reading activity by Taq DNA polymerase and the amplified DNA products contain nontemplated A overhangs at their 3' end (<http://www.invitrogen.com/>). The enzyme, buffer and MgCl₂ were included in the Taq DNA Polymerase kit (Invitrogen).

Reaction set up	Final concentration
DNA template	250 ng genomic DNA
	or 0.5-2 μ l GenomePlex library
	or 100-500 pg plasmid DNA
10xPCR buffer without MgCl ₂	1x
50 mM MgCl ₂	3 mM
10 mM dNTP (10mM each)	0.2 mM each
10 pmol/ μ l forward primer	0.2 μ M
10 pmol/ μ l reverse primer	0.2 μ M
Taq DNA polymerase (5 U/ μ l)	2 U per reaction
UltraPure water was added to the intended final volume.	

PCR system 2: PCR with Pfu DNA polymerase

The Pfu DNA polymerase has a proof-reading activity (<http://www.promega.com/>). This PCR system was used to obtain accurate sequences of the DNA targets. The products are blunt-ended. The Pfu DNA Polymerase and Pfu buffer were from Promega.

Reaction set up	Final concentration
DNA template	250 ng genomic DNA
	or 0.5-2 µl GenomePlex library
	or 100-500 pg plasmid DNA
10xPfu buffer	1x
10 mM dNTP (10mM each)	0.2 mM each
10 pmol/µl forward primer	0.2 µM
10 pmol/µl reverse primer	0.2 µM
Pfu DNA polymerase (5 U/µl)	2 U per reaction
UltraPure water was added to the intended final volume.	

PCR system 3: PCR with Expand Long Template PCR System

This system can amplify large targets up to 20 kb long. The enzyme mix consists of Taq DNA polymerase and Tgo DNA polymerase, thus possessing a proof-reading activity (<http://www.roche-applied-science.com/>). It was used for amplification of up to 8 kb DNA targets in this work, where the accuracy of sequences was considered. The products contain nontemplated A overhang at their 3' end. The enzyme mix and buffer were included in the Expand long template PCR system kit (Roche Applied Science). Reaction and cycling conditions were setup as recommended by the supplier.

Reaction set up	Final concentration
DNA template	250 ng genomic DNA
	or 0.5-2 µl GenomePlex library
	or 100-500 pg plasmid DNA
10xBuffer 3	1x
10 mM dNTP (10mM each)	0.2 mM each
10 pmol/µl forward primer	0.2 µM
10 pmol/µl reverse primer	0.2 µM
Enzyme mix	0.75 µl per reaction
UltraPure water was added to the intended final volume.	

PCR cycle program for Taq and Pfu system

1. Initial denaturation	95°C, 3 min
2. Denaturation	94°C, 30 s

- | | |
|---|-----------------------------|
| 3. Annealing | 60°C, 30 s |
| 4. Elongation | 72°C, 1 min per 1 kb target |
| 5. Repeat from step 2 for another 34 cycles | |
| 6. Final elongation | 72°C, 10 min |
| 7. Cooling | 4°C, forever |

PCR cycle program for Expand Long Template system

- | | |
|---|---------------------------------|
| 1. Initial denaturation | 94°C, 3 min |
| 2. Denaturation | 94°C, 30 s |
| 3. Annealing | 60°C, 30 s |
| 4. Elongation | 68°C, 8 min |
| 5. Repeat from step 2 for another 9 cycles | |
| 6. Denaturation | 94°C, 30 s |
| 7. Annealing | 60°C, 30 s |
| 8. Elongation | 68°C, 8 min |
| | +20 s for each successive cycle |
| 9. Repeat from step 6 for another 24 cycles | |
| 10. Final elongation | 68°C, 10 min |
| 11. Cooling | 4°C, forever |

4.15 Purification of PCR products

In some occasions, the PCR products needed to be purified to get rid of unspecific product bands or to get rid of salts in the PCR solution. This was accomplished by two alternative methods: direct PCR purification and gel extraction.

For direct PCR purification, the total PCR products were purified through columns using QIAquick PCR Purification kit according to the manufacturer's manual. This kit eliminates primer-dimer products (DNA size below 100 bp) and salts from the original PCR products.

For gel extraction, total PCR products were loaded on 1.0-1.5% TAE agarose gel and run until the band of interest was clearly separated from other bands. The target band was excised from the gel and the DNA was purified from the gel using QIAquick Gel Extraction kit according to the manufacturer's manual.

Routinely, an aliquot of the purified PCR band/product was run in agarose gel electrophoresis to ascertain its quality, before it was used for further experiments.

4.16 Cloning of PCR products

The PCR products from Taq DNA polymerase and Expand Long Template PCR systems contain nontemplated A overhangs at the 3' end, while those from Pfu DNA polymerase PCR system have blunt end. The PCR products with sticky end were cloned into TA vectors. Two commercial cloning systems for TA cloning were used: Invitrogen's TOPO TA PCR cloning kit (TopoPCR4 vector) and Stratagene's StrataClone PCR cloning kit (pSC-A vector). The blunt-end PCR products were cloned into blunt vectors with either of the two commercial kits: Invitrogen's TOPO Blunt PCR cloning kit (TopoPCR4 vector) and Stratagene's StrataClone Blunt PCR cloning kit (pSC-B vector).

The cloning procedures were performed according to the manufacturer's manual, with some modifications. For all four kits used, the cloning reaction volume was reduced to half of the recommended volume from the manufacturer. Briefly, 1 μ l PCR product was mixed with 1.5 μ l cloning buffer and 0.5 μ l vector mix in 3- μ l reaction volume, then incubated at room temperature for 15-30 minutes. The reaction was transferred into thawed competent cells (Top10 for TOPO kits, and SoloPack for StrataClone kits), mixed and incubated on ice for 30 minutes. The reaction was heat shocked at 42°C for 1 minute and then immediately put on ice. After 2 minutes on ice, 250 μ l prewarmed SOC medium were added. The reaction was incubated at 37°C, with horizontal shaking at 220 rpm, for 1 hour before spread on LB agar plates containing selective antibiotics (100 μ g/ml ampicillin). The plates were incubated overnight at 37°C to allow colonies to grow.

4.17 Plasmid DNA preparation

Quick and dirty plasmid DNA mini preparation method (quick-miniprep) was used as a rapid plasmid DNA isolation method for screening of positive clones. The plasmid DNA obtained by this method is suitable for restriction enzyme digestions, but not for direct sequencing.

To screen for positive clones, single bacterial colonies were picked from the LB plates spread after cloning. Each of them were cultured in 5 ml LB with 100 μ g/ml ampicillin at 37°C with shaking at 220 rpm overnight. Next morning, bacterial cells were collected from 1 ml of each overnight-cultures by centrifugation at 13000 rpm for 4 minutes. Cell pellets were resuspended in 100 μ l TEG buffer. To lyse the cells, 200 μ l NaOH/SDS solution was added to the cell suspension and mixed by inverting the tubes 10 times. The lysates were incubated at room temperature for 3-5 minutes until the solution became viscous. 150 μ l 3 M sodium acetate (pH 5.2) was added to each reaction, and the tubes were shortly vortexed before put on ice. After 5 minutes on ice, cell debris was separated by centrifugation at 13000 rpm for 7-10 minutes. The supernatant was

transferred to a new tube containing 800 µl ice-cold absolute ethanol and the reactions were incubated on ice for 2 minutes. The precipitated plasmid DNA in the supernatant was collected by centrifugation at 13000 rpm for 5 minutes, and then washed with 1 ml ice-cold 75% ethanol. Air-dry plasmid DNA pellets were resuspended with 30 µl TE/RNase solution. RNase was added to eliminate the co-isolated RNA. 3-µl aliquots of each plasmid DNA were digested with EcoRI, which cut at two locations on the vector adjacent to the cloning site. The digestion reactions were incubated at 37°C for 1 hour before loading on agarose gel for electrophoresis. Plasmid DNA containing correct inserts can be detected by the presence of the vector band and the insert band.

For sequencing and for long-term plasmid DNA storage, the plasmid DNAs were isolated from the bacterial cultures using commercial kits. Two formats of commercial kits were used in this study: QIAprep Spin Miniprep kit (for small scale) and Qiagen Plasmid Maxi kit (for large scale). For small scale, bacterial cell pellets were collected from 4 ml of bacterial cultures (with positive inserts). For large scale, bacterial cell pellets were collected from 100 ml of bacterial cultures (with positive inserts). The plasmid DNA isolation procedures were performed as described by the manufacturer's manuals.

4.18 DNA sequencing

Plasmid DNA containing inserts of interest were sent to Andreas Hunziker (DKFZ Genomics and Proteomics Core Facilities) for sequencing by the Sanger DNA sequencing method. The DNAs were sequenced with Big-Dye terminator chemistry on an ABI model sequencer.

4.19 Southern hybridization

Southern hybridization is a method for detection of a specific DNA sequence in a DNA population (such as genomic DNA or PCR products) by hybridizing the DNA population with a specific DNA probe (Mays Hoopes, 2008). It is divided into 4 steps: DNA transfer from agarose gel to nylon membrane (blotting), labeling of probes, hybridization and signal visualization.

Transfer of DNA from agarose gel to nylon membrane

Genomic DNA was first digested by appropriate restriction endonuclease(s) before loaded on 0.8% agarose gel. For PCR products, restriction digestion may be omitted. After separation was complete, selected DNA molecular marker positions were marked on the gel using the scalpel. The gel was soaked in 0.25 M HCl for 10 minutes to partially depurinate DNA, allowing more efficient transfer of large DNA fragments (larger than 4 kb). This step was not necessary for PCR products. Then the gel was immersed with agitation in denaturation buffer (0.4 M NaOH) twice,

each for 15 minutes, to denature the DNA. Afterward the gel was ready for transfer. Vertical downward capillary transfer was used in this study.

In preparation for transfer, one Hybond N+ nylon membrane was cut to the size slightly larger than the gel and four Whatman 3 MM papers were cut to the size slightly larger than the nylon membrane. The system was setup by placing a stack of tissue paper on the table first, followed by two dry Whatman paper. Another two Whatman papers were soaked in 0.4 M NaOH before placed on top of the dry Whatman papers. The nylon membrane was shortly soaked in 0.4 M NaOH and then placed on top of the Whatman papers. The gel from previous denaturation step was placed atop the nylon membrane. For efficient transfer, it is important to avoid any air bubbles during the transfer setup. The whole setup was wrapped with plastic foil. A small plastic board (only large enough to cover the gel) was placed on top of the setup and then 250 g weight was placed on top of the board. The DNA transferring process took approximately 4-6 hours. During the transfer, DNA was driven from the gel to the membrane by capillary force. Once arrived at the membrane, the DNA binds to the positively charged nylon membrane, and is fixed to the membrane by the alkaline solvent. After the transfer was complete, the marks of the selected DNA markers, as well as the wells on the gel, were marked on the nylon membrane with colored pencil. The membrane was then soaked in the neutralizing buffer (3 M NaCl, 0.3 M sodium citrate, 0.5 M Tris-HCl, pH 7.0) for 15 minutes before let dry on a clean Whatman paper at room temperature.

Radioactive labeling of DNA probes

For hybridization, the DNA probes were labeled with [α - 32 P] dCTP, using Rediprime II DNA Labeling System (GE Healthcare) as described in the manufacturer's manual. Briefly, 25-30 ng linear DNA probe was diluted in UltraPure water to a volume of 45 μ l. The DNA was denatured at 95°C for 5 minutes, then put immediately on ice. The DNA was spun down shortly, and transferred to a Rediprime II reaction mixture tube. 5 μ l of 5 μ Ci [α - 32 P] dCTP was added to the mixture and mixed by pipetting. The reaction was incubated at 37°C for 30 minutes and stopped by adding 20 μ l stop-solution (2xTNE, 0.1% BPB, 0.5% blue dextran). The labeling reactions were later purified through pre-equilibrated columns. To equilibrate a purification column, the resin particles in a MicroSpin G-25 column were resuspended by vigorous vortex. The tip of the column was cut off, and the column was placed in a 2 ml collecting tube. The tubes were centrifuged at 300 rpm for 1 minute and the flow-through was discarded. To equilibrate the resin, 200 μ l 1xTNE was added to the column, and the tubes were centrifuged at 300 rpm for 1 minute. The column was transferred into a new collecting tube. 75 μ l 1xTNE were added to the complete labeling reaction. The solution was mixed by pipetting and transferred to the equilibrated column. The tubes were centrifuged at 300 rpm for 1 minute. The small DNA and dNTP remained in the

matrix of the column, while larger DNA molecules passed through. The eluate was stored at -20°C or directly used for hybridization.

DNA hybridization with a radioactive labeled probe

The hybridization of labeled DNA probe with the complementary sequences on the membrane were performed with high-stringency condition, at 68°C. The hybridization oven and hybridization solution were prewarmed to 68°C before use. The membrane was prehybridized at 68°C in the oven with constant rotation for at least 1 hour before it was ready for hybridization. Prior to adding the labeled probe, the prehybridization solution was discarded and 10 µl of fresh prewarmed hybridization solution was added to the tube. 70 µl of the labeled probe were used per 10 ml hybridization volume. The labeled probe was denatured by the addition of 0.2 M NaOH (final concentration) and incubated at room temperature for 10 minutes. The denatured probe was added to the hybridization tube, and the membrane was hybridized overnight at 68°C with constant rotation. Next day, the membrane was washed 3x 30 minutes with prewarmed washing buffer (2xSSC, 0.1% SDS) at 68°C with constant rotation to reduce unspecific binding background. Afterward, the membrane was dried at room temperature and wrapped with plastic foil.

Detection of hybridization signal by autoradiography

The P-32 emits high-energy β-particles that can easily be detected by autoradiography. To set up the exposure, the membrane was taped to a piece of Whatman 3MM paper and placed in a cassette. An autoradiography film, MR film, was placed over the membrane and covered with an intensifying screen. The cassette was closed and stored at -80°C. Genomic DNA hybridization requires longer exposure time than PCR product hybridization. The usual exposition times were 7 days and 30-120 minutes for genomic DNA and PCR products, respectively. After the exposure was complete, the film was taken out and developed by Curix 60 film processor machine.

4.20 DNA amplification with phi29 DNA polymerase

Phi29 DNA polymerase is a replicative polymerase isolated from *Bacillus subtilis* phage phi29. The enzyme exhibits exceptional strand-displacement and processive synthesis properties. It has also inherent 3'-5' proofreading exonuclease activity (Blanco et al, 1989; Esteban et al, 1993). The enzyme is able to amplify from both circular and linear DNA templates at 30°C, and the primers are required to be exonuclease-resistant (Dean et al, 2002; Hutchison et al, 2005).

In this work, phi29 DNA polymerase, in combination with random hexamer, was used to amplified linear DNA template (called multiple displacement amplification, MDA) and circular

DNA template (called rolling circle amplification, RCA). To make the random hexamer resistant to exonuclease activity, it was modified with phosphothioate at the 3' end (Hutchison et al, 2005). The reaction conditions of both MDA and RCA were identical except for the nature of template DNA. The products of MDA and RCA are high molecular weight linear DNA and high molecular weight linear concatameric DNA respectively. The DNA may reach up to 70 kb long.

The reactions were carried out at 20 µl final volume. The maximum volume of template DNA was limited to 3 µl. The volume of the template DNA was brought to 3 µl by adding UltraPure water. First, the template DNA was mixed with 1 µl 500 pmol/µl random hexamer (phosphothioate modified at 3' end) and 1 µl 5xAnnealing buffer. The reaction was mixed by pipetting and denatured at 95°C for 3 minutes and immediately left on ice for 2 minutes. The polymerase started by adding 2 µl 10xBSA, 2 µl 10xPhi29 buffer, 5 µl 10 mM dNTP, 1 µl phi29 DNA polymerase (Fermentas) and 5 µl UltraPure water. The reaction was incubated at 30°C for 16 hours, then inactivated at 65°C for 10 minutes. The quality of the amplified DNA products was observed by running a 1-µl aliquot on 1% agarose gel together with lamda-DNA/HindIII marker. The products should distribute mostly between 23 kb band and the 2 kb band of the DNA marker.

4.21 Preparation of pBluescript (pBS) vector for cloning at EcoRI

To prepare the pBS vector for cloning at the EcoRI site, 12 µg pBS was digested with 140 U EcoRI in 100 µl reaction volume. The linearized vector was purified by phenol/chloroform-chloroform extraction and later precipitated with sodium acetate and ethanol. In brief, 1 volume of phenol/chloroform solution was added to the digestion reaction. The mixture was vigorously shaken for 15 seconds and centrifuged at 13000 for 5 minutes. The upper aqueous layer was collected in a new tube, and 1 volume of chloroform was added. The reaction tube was vigorously shaken for 15 seconds and centrifuged at 13000 rpm for 5 minutes. The upper aqueous phase was collected in a new tube and the plasmid DNA was precipitated by addition of 0.1 volume 3 M sodium acetate (pH 5.2) and 2.5 volumes of ice-cold absolute ethanol. The tube was incubated at -20°C for 1 hour before centrifuged at 13000 rpm for 20 minutes. The DNA pellet was washed with 300 µl ice-cold 75% ethanol. The air-dried DNA pellet was resuspended in UltraPure water.

The linearized pBS was dephosphorylated to prevent self-ligation. Briefly, a 4 µg aliquot of the linearized pBS was added with 5 µl 10xCiAP buffer, 5 µl CiAP (1 U/µl) and UltraPure water to final volume of 50 µl. The reaction was incubated at 37°C for 2 hours before inactivated at 85°C for 20 minutes. The reaction was purified by phenol/chloroform-chloroform extraction and later precipitated with sodium acetate and ethanol, as described above. The dephosphorylated EcoRI-

linearized pBS vector was resuspended with 12 µl UltraPure water and stored at -20°C prior to use.

4.22 ASP16 strategy for HPV16 integration analysis

The newly developed ASP16 strategy for HPV16 integration analysis in cervical scrapes (Xu, 2010) consists of four major steps: (1) whole genome amplification, (2) HPV16 DNA enrichment, (3) high-throughput DNA sequencing using the Roche/454 Genome Sequencer FLX, and (4) data analysis.

4.22.1 GenomePlex whole genome amplification

Since DNA of clinical samples was only available in nanogram amounts, the genomic DNA samples were amplified with GenomePlex whole genome amplification kit (GenomePlex WGA2 kit, Sigma Aldrich) to provide a library of amplified DNA fragments. In the GenomePlex WGA method, the DNA is first chemically fragmented, then GenomePlex universal adapters are added to both ends of the fragmented DNA molecules to create template molecules for library amplification, and the complete library is obtained by PCR amplification of these template molecules. Such library is called OmniPlex library. As the result of fragmentation, the OmniPlex library contained amplified DNA molecules with size distribution between 200 bp - 1.5 kb. The GenomePlex amplification was carried out according to the manufacturer's manual. The minimum template DNA was 10 ng. Generally for clinical DNA samples, 10-50 ng DNA were used for templates. The entire incubation process was carried out on a Thermocycler PTC-200 machine.

4.22.2 HPV16 enrichment from GenomePlex DNA libraries

The enrichment of HPV16 DNA was carried out in three steps: linear amplification of HPV16 DNA, isolation of HPV16 DNA via biotin-streptavidin, and HPV16 semi-nested PCR. In this study, the HPV16-specific linear amplification and semi-nested PCR were performed as multiplex amplification reactions. Four amplification sets were performed for each DNA sample during both of these steps. Each amplification set consisted of a combination of 7-8 HPV16 primers. Primers are listed in section 4.7.4 (Table 4.4, Table 4.5 and Table 4.6).

HPV16 linear amplification by primer extension

Thirty-two 5'-biotin-labeled HPV16-specific forward primers, located in the E1-E2 region (see Figure 2.33), were used in this step, dividing into four groups of primer mixes. See Table 2.14 in

Results section 2.2.2.3 for primer combinations of the four mixes. Four linear amplification reactions were setup for each DNA sample, each reaction containing one of the four primer mixes, respectively. In the primer mixes, each biotin-labeled HPV16 primer had a concentration of 0.1 μ M. Linear amplification was carried out in 50 μ l reaction volume, using Taq DNA polymerase (Invitrogen). In brief, 2.5 μ l GenomePlex library was mixed with 5 μ l 10xreaction buffer without $MgCl_2$, 3 μ l 50mM $MgCl_2$, 1 μ l dNTP (10 mM each), 8 μ l biotin-labeled HPV16 primer mix, 0.5 μ l Taq DNA polymerase and 30 μ l UltraPure water. The polymerization reaction was carried out for 45 cycles. The amplification products were single-stranded HPV16 DNA with biotin-labeled at the 5' end of the DNA strand. The cycle conditions are described below:

- | | |
|---|--------------|
| 1. Initial denaturation | 95°C, 2 min |
| 2. Denaturation | 94°C, 30 s |
| 3. Annealing | 60°C, 75 s |
| 4. Elongation | 72°C, 40 s |
| 5. Repeat steps 2-4 for another 44 cycles | |
| 6. Cooling | 4°C, forever |

Isolation of HPV16 DNA with biotin-streptavidin system

The biotin at the 5' end of each HPV16 single-stranded DNA molecules allowed easy isolation of these molecules with streptavidin-coated magnetic beads. BioMag streptavidin beads (Polysciences) were used. Briefly, a 75- μ l aliquot of suspended BioMag streptavidin beads (1 mg/ml) was transferred to a 1.5-ml tube. The tube was placed in a magnet and incubated for 1 minute, allowing the beads to concentrate at one side of the tube. The supernatant was discarded and the tube was removed from the magnet. To wash the beads, they were resuspended in 100 μ l 2xBinding buffer (40 mM Tris-HCl, 2 mM EDTA, 2 M NaCl, 0.1% Tween20, pH 7.8 in DEPC water), separated from supernatant on the magnet, and resuspended again in 50 μ l 2xBinding buffer. A total 50 μ l reaction from the linear amplification step was added to the bead suspension, mixed by pipetting, and incubated for 15-30 minutes at room temperature with occasional agitation. The biotin-labeled HPV16 DNA molecules bound to the beads, while the background genomic DNA did not. The HPV16-DNA-bound BioMag beads were placed on the magnet, incubated for 1 minute and the supernatant was discarded. The beads were washed twice with 150 μ l 1xWash buffer (20 mM Tris-HCl, 1 mM EDTA, 1 M NaCl, 0.05% Tween20, pH 7.8 in DEPC water), then three times with 150 μ l 1xTE buffer (in DEPC water). The beads were resuspended in 30 μ l 1xTE buffer and stored at -20°C.

HPV16 semi-nested multiplex PCR

The isolated single-stranded HPV16 DNA molecules were amplified by multiplex PCR, generating amplified double-stranded HPV16 DNA as the products. To make the PCR products

compatible for GS-FLX sequencing, Roche-A and Roche-B adapters were required at both ends of each DNA molecule. As forward primers, nested HPV16 primers, fused with the Roche-A sequence at the 5' end, were used (RA_HP16, see Table 4.5). The reverse primers contained the Roche-B sequence at the 5' end, followed by 4-nt barcode and then the GenomePlex universal adapter sequence (RB-B01 to -B24, see Table 4.6). For each DNA sample, the nested HPV16 primers were assembled in four groups and used in four separate multiplex PCR reactions. For each DNA sample, one reverse primer with a unique barcode was used. See Table 2.14 in Results section 2.2.2.3 for primer combinations of the four reactions. The multiplex PCR was carried out in 50 µl reaction volume, using Qiagen multiplex PCR kit, according to the manufacturer's manual. Briefly, a 2.5 µl aliquot of each purified single-stranded HPV16 DNA product from previous step was mixed with 25 µl 2xQiagen multiplex PCR master mix, 2.5 µl of corresponding 20xRA_HP16 primer mix (2 µM each), 1 µl 10 pmol/µl RB-Barcode-GPUA primer (unique for each DNA sample) and 19 µl UltraPure water. The cycle conditions are described below:

- | | |
|---|--------------|
| 1. Initial denaturation | 95°C, 15 min |
| 2. Denaturation | 94°C, 30 s |
| 3. Annealing | 60°C, 75 s |
| 4. Elongation | 72°C, 1 min |
| 5. Repeat from step 2 for another 34 cycles | |
| 6. Cooling | 4°C, forever |

The quality of the multiplex PCR products was determined by running 3-µl aliquots on agarose gels. By blotting the gel and hybridizing with an HPV16-specific probe, the specificity of the multiplex PCR products was determined.

4.22.3 Size selection of HPV16 multiplex PCR products

In experiment ASP16-4, it was attempted to limit the size distribution of the multiplex PCR product to 200-300 bp, using the E-Gel SizeSelect system (Invitrogen). The procedure was carried out according to the manufacturer's manual. Each multiplex PCR product was loaded on a single well of E-Gel SizeSelect. The E-Gel contained 2% agarose with dye that allowed real-time visualization of DNA migration. Smear portions of the desired DNA size were collected from the collecting wells by directly pipetting the solution out. The quality and specificity of the smear portion of each multiplex PCR products were determined by agarose gel electrophoresis and Southern hybridization with HPV16 probe. The four size-selected aliquots of each DNA sample were pooled in equal amounts, approximately 100 ng each. The mixture was concentrated by precipitation with 3 M sodium acetate (pH 5.2) and ethanol, and dissolved in 100 µl TE. The concentration and quality of the pooled amplicons was measured with NanoDrop, before it was sent to DKFZ Genomics and Proteomics Core Facilities for sequencing.

4.22.4 Roche/454 GS-FLX pyrosequencing

The amplicons were sequenced by the Roche/454 GS-FLX pyrosequencing. One of the two regions on 70x75 picotiter plate was chosen for sequencing platform. According to the manufacturer, the selected sequencing format can generate, per each sequencing round, up to 210000 sequence reads of up to 250 nt in length. Individual amplicon molecules were amplified by emulsion PCR (emPCR) before they were sequenced on the sequencing platform (more details can be found at <http://www.454.com/>). In this work, Roche-B was used as the sequencing primer. This resulted in sequence reads starting from the 4-nt barcode sequence, followed by GenomePlex universal adapter, then the amplicon sequence, and then the nested RA-HPV16 primer sequence at the end. The outputs were generated as a list of FASTA sequence reads.

4.22.5 Data analysis

Each Roche/454 GS-FLX run generated up to 210000 single sequence reads. This required initial automatic data analysis and manipulation. Writing a suitable set of programs to analyze ASP16 sequence data using Perl programming language was an essential part of this PhD work. The details are presented in Results and Discussion sections. Shortly, the output FASTA sequence reads from GS-FLX sequencing were sorted into sample groups according to the 4-nt barcodes. Then for each sample group, the programs analyzed and selected only informative sequence reads into output lists. These selected sequences were aligned with Clustal algorithms, and then re-edited for easier visualization. The sequence and sequence alignment outputs can be viewed and manipulated with any sequence editor software. In this work, HUSAR and Geneious Pro were employed. The programs then proceeded to calculate and generate basic statistical data, of which the outputs were written in tab-delimited format that may be open in Microsoft Excel or similar software.

References

- Azizi N, Brazete J, Hankins C, Money D, Fontaine J, Koushik A, Rachlis A, Pourreaux K, Ferenczy A, Franco E, Coutlee F (2008) Influence of human papillomavirus type 16 (HPV-16) E2 polymorphism on quantification of HPV-16 episomal and integrated DNA in cervicovaginal lavages from women with cervical intraepithelial neoplasia. *The Journal of general virology* **89**: 1716-1728
- Baker CC, Phelps WC, Lindgren V, Braun MJ, Gonda MA, Howley PM (1987) Structural and transcriptional analysis of human papillomavirus type 16 sequences in cervical carcinoma cell lines. *Journal of virology* **61**: 962-971
- Batta K, Kundu TK (2007) Activation of p53 function by human transcriptional coactivator PC4: role of protein-protein interaction, DNA bending, and posttranslational modifications. *Molecular and cellular biology* **27**: 7603-7614
- Batta K, Yokokawa M, Takeyasu K, Kundu TK (2009) Human transcriptional coactivator PC4 stimulates DNA end joining and activates DSB repair activity. *Journal of molecular biology* **385**: 788-799
- Bernard HU, Burk RD, Chen Z, van Doorslaer K, Hausen H, de Villiers EM (2010) Classification of papillomaviruses (PVs) based on 189 PV types and proposal of taxonomic amendments. *Virology* **401**: 70-79
- Blanco L, Bernad A, Lazaro JM, Martin G, Garmendia C, Salas M (1989) Highly efficient DNA synthesis by the phage phi 29 DNA polymerase. Symmetrical mode of DNA replication. *The Journal of biological chemistry* **264**: 8935-8940
- Bory JP, Cucherousset J, Lorenzato M, Gabriel R, Quereux C, Birembaut P, Clavel C (2002) Recurrent human papillomavirus infection detected with the hybrid capture II assay selects women with normal cervical smears at risk for developing high grade cervical lesions: a longitudinal study of 3,091 women. *International journal of cancer* **102**: 519-525
- Boyer SN, Wazer DE, Band V (1996) E7 protein of human papilloma virus-16 induces degradation of retinoblastoma protein through the ubiquitin-proteasome pathway. *Cancer research* **56**: 4620-4624
- Briolat J, Dalstein V, Saunier M, Joseph K, Caudroy S, Pretet JL, Birembaut P, Clavel C (2007) HPV prevalence, viral load and physical state of HPV-16 in cervical smears of patients with different grades of CIN. *International journal of cancer* **121**: 2198-2204
- Callahan DE, Karim A, Zheng G, Tso PO, Lesko SA (1992) Quantitation and mapping of integrated human papillomavirus on human metaphase chromosomes using a fluorescence microscope imaging system. *Cytometry* **13**: 453-461
- Calleja-Macias IE, Kalantari M, Allan B, Williamson AL, Chung LP, Collins RJ, Zuna RE, Dunn ST, Ortiz Lopez R, Barrera-Saldana HA, Cubie HA, Cuschieri K, Villa LL, Bernard HU (2005) Papillomavirus subtypes are natural and old taxa: phylogeny of human papillomavirus types 44 and 55 and 68a and -b. *Journal of virology* **79**: 6565-6569
- Casas L, Galvan SC, Ordonez RM, Lopez N, Guido M, Berumen J (1999) Asian-american variants of human papillomavirus type 16 have extensive mutations in the E2 gene and are highly amplified in cervical carcinomas. *International journal of cancer* **83**: 449-455
- Castellsague X (2008) Natural history and epidemiology of HPV infection and cervical cancer. *Gynecologic oncology* **110**: S4-7
- Chan JK, Monk BJ, Brewer C, Keefe KA, Osann K, McMeekin S, Rose GS, Youssef M, Wilczynski SP, Meyskens FL, Berman ML (2003) HPV infection and number of lifetime sexual partners are strong predictors for 'natural' regression of CIN 2 and 3. *British journal of cancer* **89**: 1062-1066

Chan SY, Ho L, Ong CK, Chow V, Drescher B, Durst M, ter Meulen J, Villa L, Luande J, Mgaya HN, et al. (1992) Molecular variants of human papillomavirus type 16 from four continents suggest ancient pandemic spread of the virus and its coevolution with humankind. *Journal of virology* **66**: 2057-2066

Chellappan S, Kraus VB, Kroger B, Munger K, Howley PM, Phelps WC, Nevins JR (1992) Adenovirus E1A, simian virus 40 tumor antigen, and human papillomavirus E7 protein share the capacity to disrupt the interaction between transcription factor E2F and the retinoblastoma gene product. *Proceedings of the National Academy of Sciences of the United States of America* **89**: 4549-4553

Chen Z, Terai M, Fu L, Herrero R, DeSalle R, Burk RD (2005) Diversifying selection in human papillomavirus type 16 lineages based on complete genome analyses. *Journal of virology* **79**: 7014-7023

Clavel C, Masure M, Bory JP, Putaud I, Mangeonjean C, Lorenzato M, Nazeyrollas P, Gabriel R, Quereux C, Birembaut P (2001) Human papillomavirus testing in primary screening for the detection of high-grade cervical lesions: a study of 7932 women. *British journal of cancer* **84**: 1616-1623

Couturier J, Sastre-Garau X, Schneider-Maunoury S, Labib A, Orth G (1991) Integration of papillomavirus DNA near myc genes in genital carcinomas and its consequences for proto-oncogene expression. *Journal of virology* **65**: 4534-4538

Coward E, Haas SA, Vingron M (2002) SpliceNest: visualization of gene structure and alternative splicing based on EST clusters. *Trends Genet* **18**: 53-55

Cullen AP, Reid R, Campion M, Lorincz AT (1991) Analysis of the physical state of different human papillomavirus DNAs in intraepithelial and invasive cervical neoplasm. *Journal of virology* **65**: 606-612

de Sanjose S, Quint WG, Alemany L, Geraets DT, Klaustermeier JE, Lloveras B, Tous S, Felix A, Bravo LE, Shin HR, Vallejos CS, de Ruiz PA, Lima MA, Guimera N, Clavero O, Alejo M, Llombart-Bosch A, Cheng-Yang C, Tatti SA, Kasamatsu E, Iljazovic E, Odida M, Prado R, Seoud M, Grce M, Usubutun A, Jain A, Suarez GA, Lombardi LE, Banjo A, Menendez C, Domingo EJ, Velasco J, Nessa A, Chichareon SC, Qiao YL, Lerma E, Garland SM, Sasagawa T, Ferrera A, Hammouda D, Mariani L, Pelayo A, Steiner I, Oliva E, Meijer CJ, Al-Jassar WF, Cruz E, Wright TC, Puras A, Llave CL, Tzardi M, Agorastos T, Garcia-Barriola V, Clavel C, Ordi J, Andujar M, Castellsague X, Sanchez GI, Nowakowski AM, Bornstein J, Munoz N, Bosch FX (2010) Human papillomavirus genotype attribution in invasive cervical cancer: a retrospective cross-sectional worldwide study. *The lancet oncology* **11**: 1048-1056

de Villiers EM, Fauquet C, Broker TR, Bernard HU, zur Hausen H (2004) Classification of papillomaviruses. *Virology* **324**: 17-27

Dean FB, Hosono S, Fang L, Wu X, Faruqi AF, Bray-Ward P, Sun Z, Zong Q, Du Y, Du J, Driscoll M, Song W, Kingsmore SF, Egholm M, Lasken RS (2002) Comprehensive human genome amplification using multiple displacement amplification. *Proceedings of the National Academy of Sciences of the United States of America* **99**: 5261-5266

Doorbar J (2005) The papillomavirus life cycle. *J Clin Virol* **32 Suppl 1**: S7-15

Doorbar J (2006) Molecular biology of human papillomavirus infection and cervical cancer. *Clin Sci (Lond)* **110**: 525-541

Doorbar J, Ely S, Sterling J, McLean C, Crawford L (1991) Specific interaction between HPV-16 E1-E4 and cytokeratins results in collapse of the epithelial cell intermediate filament network. *Nature* **352**: 824-827

Durst M, Croce CM, Gissmann L, Schwarz E, Huebner K (1987) Papillomavirus sequences integrate near cellular oncogenes in some cervical carcinomas. *Proceedings of the National Academy of Sciences of the United States of America* **84**: 1070-1074

Eriksson A, Herron JR, Yamada T, Wheeler CM (1999) Human papillomavirus type 16 variant lineages characterized by nucleotide sequence analysis of the E5 coding segment and the E2 hinge region. *The Journal of general virology* **80 (Pt 3)**: 595-600

- Esteban JA, Salas M, Blanco L (1993) Fidelity of phi 29 DNA polymerase. Comparison between protein-primed initiation and DNA polymerization. *The Journal of biological chemistry* **268**: 2719-2726
- Fehrmann F, Klumpp DJ, Laimins LA (2003) Human papillomavirus type 31 E5 protein supports cell cycle progression and activates late viral functions upon epithelial differentiation. *Journal of virology* **77**: 2819-2831
- Ferber MJ, Thorland EC, Brink AA, Rapp AK, Phillips LA, McGovern R, Gostout BS, Cheung TH, Chung TK, Fu WY, Smith DI (2003) Preferential integration of human papillomavirus type 18 near the ϵ -myc locus in cervical carcinoma. *Oncogene* **22**: 7233-7242
- Fu L, Teraï M, Matsukura T, Herrero R, Burk RD (2004) Codetection of a mixed population of candHPV62 containing wild-type and disrupted E1 open-reading frame in a 45-year-old woman with normal cytology. *The Journal of infectious diseases* **190**: 1303-1309
- Giannoudis A, Herrington CS (2001) Human papillomavirus variants and squamous neoplasia of the cervix. *The Journal of pathology* **193**: 295-302
- Graham DA, Herrington CS (2000) HPV-16 E2 gene disruption and sequence variation in CIN 3 lesions and invasive squamous cell carcinomas of the cervix: relation to numerical chromosome abnormalities. *Mol Pathol* **53**: 201-206
- Ham J, Dostatni N, Gauthier JM, Yaniv M (1991) The papillomavirus E2 protein: a factor with many talents. *Trends in biochemical sciences* **16**: 440-444
- Ho GY, Bierman R, Beardsley L, Chang CJ, Burk RD (1998) Natural history of cervicovaginal papillomavirus infection in young women. *The New England journal of medicine* **338**: 423-428
- Ho L, Chan SY, Burk RD, Das BC, Fujinaga K, Icenogle JP, Kahn T, Kiviat N, Lancaster W, Mavromara-Nazos P, et al. (1993) The genetic drift of human papillomavirus type 16 is a means of reconstructing prehistoric viral spread and the movement of ancient human populations. *Journal of virology* **67**: 6413-6423
- Hopman AH, Smedts F, Dignef W, Ummelen M, Sonke G, Mravunac M, Vooijs GP, Speel EJ, Ramaekers FC (2004) Transition of high-grade cervical intraepithelial neoplasia to micro-invasive carcinoma is characterized by integration of HPV 16/18 and numerical chromosome abnormalities. *The Journal of pathology* **202**: 23-33
- Hudelist G, Manavi M, Pischinger KI, Watkins-Riedel T, Singer CF, Kubista E, Czerwenka KF (2004) Physical state and expression of HPV DNA in benign and dysplastic cervical tissue: different levels of viral integration are correlated with lesion grade. *Gynecologic oncology* **92**: 873-880
- Hutchison CA, 3rd, Smith HO, Pfannkoch C, Venter JC (2005) Cell-free cloning using phi29 DNA polymerase. *Proceedings of the National Academy of Sciences of the United States of America* **102**: 17332-17336
- Jacobs MV, de Roda Husman AM, van den Brule AJ, Snijders PJ, Meijer CJ, Walboomers JM (1995) Group-specific differentiation between high- and low-risk human papillomavirus genotypes by general primer-mediated PCR and two cocktails of oligonucleotide probes. *Journal of clinical microbiology* **33**: 901-905
- Jansen R, Yu H, Greenbaum D, Kluger Y, Krogan NJ, Chung S, Emili A, Snyder M, Greenblatt JF, Gerstein M (2003) A Bayesian networks approach for predicting protein-protein interactions from genomic data. *Science* **302**: 449-453
- Jarrett CR, Blancato J, Cao T, Bressette DS, Cepeda M, Young PE, King CR, Byers SW (2001) Human APC2 localization and allelic imbalance. *Cancer research* **61**: 7978-7984
- Jeon S, Allen-Hoffmann BL, Lambert PF (1995) Integration of human papillomavirus type 16 into the human genome correlates with a selective growth advantage of cells. *Journal of virology* **69**: 2989-2997

- Jeon S, Lambert PF (1995) Integration of human papillomavirus type 16 DNA into the human genome leads to increased stability of E6 and E7 mRNAs: implications for cervical carcinogenesis. *Proceedings of the National Academy of Sciences of the United States of America* **92**: 1654-1658
- Jiang M, Baseman JG, Koutsky LA, Feng Q, Mao C, Kiviat NB, Xi LF (2009) Sequence variation of human papillomavirus type 16 and measurement of viral integration by quantitative PCR. *Journal of clinical microbiology* **47**: 521-526
- Jiang X, Hanna Z, Kaouass M, Girard L, Jolicoeur P (2002) Ahi-1, a novel gene encoding a modular protein with WD40-repeat and SH3 domains, is targeted by the Ahi-1 and Mis-2 provirus integrations. *Journal of virology* **76**: 9046-9059
- Kalantari M, Blennow E, Hagmar B, Johansson B (2001) Physical state of HPV16 and chromosomal mapping of the integrated form in cervical carcinomas. *Diagn Mol Pathol* **10**: 46-54
- Kammer C, Tommasino M, Syrjanen S, Delius H, Hebling U, Warthorst U, Pfister H, Zehbe I (2002) Variants of the long control region and the E6 oncogene in European human papillomavirus type 16 isolates: implications for cervical disease. *British journal of cancer* **86**: 269-273
- Kammer C, Warthorst U, Torrez-Martinez N, Wheeler CM, Pfister H (2000) Sequence analysis of the long control region of human papillomavirus type 16 variants and functional consequences for P97 promoter activity. *The Journal of general virology* **81**: 1975-1981
- Kimbauer R, Taub J, Greenstone H, Roden R, Durst M, Gissmann L, Lowy DR, Schiller JT (1993) Efficient self-assembly of human papillomavirus type 16 L1 and L1-L2 into virus-like particles. *Journal of virology* **67**: 6929-6936
- Klaes R, Woerner SM, Ridder R, Wentzensen N, Duerst M, Schneider A, Lotz B, Melsheimer P, von Knebel Doeberitz M (1999) Detection of high-risk cervical intraepithelial neoplasia and cervical cancer by amplification of transcripts derived from integrated papillomavirus oncogenes. *Cancer research* **59**: 6132-6136
- Kraus I, Driesch C, Vinokurova S, Hovig E, Schneider A, von Knebel Doeberitz M, Durst M (2008) The majority of viral-cellular fusion transcripts in cervical carcinomas cotranscribe cellular sequences of known or predicted genes. *Cancer research* **68**: 2514-2522
- Kuslich CD, Chui B, Yamashiro CT (2008) Overview of PCR. *Curr Protoc Essential Lab Tech* **00**: 10.12.11-31
- Li N, Franceschi S, Howell-Jones R, Snijders PJ, Clifford GM (2010) Human papillomavirus type distribution in 30,848 invasive cervical cancers worldwide: Variation by geographical region, histological type and year of publication. *International journal of cancer*
- Li W, Wang W, Si M, Han L, Gao Q, Luo A, Li Y, Lu Y, Wang S, Ma D (2008) The physical state of HPV16 infection and its clinical significance in cancer precursor lesion and cervical carcinoma. *Journal of cancer research and clinical oncology* **134**: 1355-1361
- Longuet M, Beaudenon S, Orth G (1996) Two novel genital human papillomavirus (HPV) types, HPV68 and HPV70, related to the potentially oncogenic HPV39. *Journal of clinical microbiology* **34**: 738-744
- Lorenzato M, Bory JP, Cucherousset J, Nou JM, Bouttens D, Thil C, Dez F, Evrard G, Quereux C, Birembaut P, Clavel C (2002) Usefulness of DNA ploidy measurement on liquid-based smears showing conflicting results between cytology and high-risk human papillomavirus typing. *American journal of clinical pathology* **118**: 708-713
- Lorenzato M, Clavel C, Masure M, Nou JM, Bouttens D, Evrard G, Bory JP, Maugard B, Quereux C, Birembaut P (2001) DNA image cytometry and human papillomavirus (HPV) detection help to select smears at high risk of high-grade cervical lesions. *The Journal of pathology* **194**: 171-176

- Luft F, Klaes R, Nees M, Durst M, Heilmann V, Melsheimer P, von Knebel Doeberitz M (2001) Detection of integrated papillomavirus sequences by ligation-mediated PCR (DIPS-PCR) and molecular characterization in cervical cancer cells. *International journal of cancer* **92**: 9-17
- Manchester KM, Heston WD, Donner DB (1993) Tumour necrosis factor-induced cytotoxicity is accompanied by intracellular mitogenic signals in ME-180 human cervical carcinoma cells. *The Biochemical journal* **290** (Pt 1): 185-190
- Matsumoto K, Yoshikawa H, Nakagawa S, Tang X, Yasugi T, Kawana K, Sekiya S, Hirai Y, Kukimoto I, Kanda T, Taketani Y (2000) Enhanced oncogenicity of human papillomavirus type 16 (HPV16) variants in Japanese population. *Cancer letters* **156**: 159-165
- Maufort JP, Shai A, Pitot HC, Lambert PF (2010) A role for HPV16 E5 in cervical carcinogenesis. *Cancer research* **70**: 2924-2931
- Mays Hoopes LL (2008) Nucleic acid blotting: Southern and Northern. *Curr Protoc Essential Lab Tech* **00**: 8.2.1-24
- McBride AA (2008) Replication and partitioning of papillomavirus genomes. *Advances in virus research* **72**: 155-205
- McBride AA, Oliveira JG, McPhillips MG (2006) Partitioning viral genomes in mitosis: same idea, different targets. *Cell cycle* **5**: 1499-1502
- McIntosh PB, Laskey P, Sullivan K, Davy C, Wang Q, Jackson DJ, Griffin HM, Doorbar J (2010) E1--E4-mediated keratin phosphorylation and ubiquitylation: a mechanism for keratin depletion in HPV16-infected epithelium. *Journal of cell science* **123**: 2810-2822
- Meissner JD (1999) Nucleotide sequences and further characterization of human papillomavirus DNA present in the CaSki, SiHa and HeLa cervical carcinoma cell lines. *The Journal of general virology* **80** (Pt 7): 1725-1733
- Metzker ML (2010) Sequencing technologies - the next generation. *Nature reviews genetics* **11**: 31-46
- Mincheva A, Gissmann L, zur Hausen H (1987) Chromosomal integration sites of human papillomavirus DNA in three cervical cancer cell lines mapped by in situ hybridization. *Medical microbiology and immunology* **176**: 245-256
- Moody CA, Laimins LA (2010) Human papillomavirus oncoproteins: pathways to transformation. *Nature reviews cancer* **10**: 550-560
- Munemitsu S, Albert I, Souza B, Rubinfeld B, Polakis P (1995) Regulation of intracellular beta-catenin levels by the adenomatous polyposis coli (APC) tumor-suppressor protein. *Proceedings of the National Academy of Sciences of the United States of America* **92**: 3046-3050
- Munger K, Phelps WC, Bubb V, Howley PM, Schlegel R (1989) The E6 and E7 genes of the human papillomavirus type 16 together are necessary and sufficient for transformation of primary human keratinocytes. *Journal of virology* **63**: 4417-4421
- Munoz N, Bosch FX, de Sanjose S, Herrero R, Castellsague X, Shah KV, Snijders PJ, Meijer CJ (2003) Epidemiologic classification of human papillomavirus types associated with cervical cancer *The New England journal of medicine* **348**: 518-527
- Nasiell K, Roger V, Nasiell M (1986) Behavior of mild cervical dysplasia during long-term follow-up. *Obstetrics and gynecology* **67**: 665-669
- Nishikawa K, Rosenblum MG, Newman RA, Pandita TK, Hittelman WN, Donato NJ (1992) Resistance of human cervical carcinoma cells to tumor necrosis factor correlates with their increased sensitivity to cisplatin: evidence of a role for DNA repair and epidermal growth factor receptor. *Cancer research* **52**: 4758-4765

- Nurnberg W, Artuc M, Nawrath M, Lovric J, Stuting S, Moelling K, Czarnetzki BM, Schadendorf D (1995) Human c-myc is expressed in cervical carcinomas and transactivates the HPV-16 promoter. *Cancer research* **55**: 4432-4437
- Ordóñez RM, Espinosa AM, Sánchez-González DJ, Armendariz-Borunda J, Berumen J (2004) Enhanced oncogenicity of Asian-American human papillomavirus 16 is associated with impaired E2 repression of E6/E7 oncogene transcription. *The Journal of general virology* **85**: 1433-1444
- Park JS, Hwang ES, Park SN, Ahn HK, Um SJ, Kim CJ, Kim SJ, Namkoong SE (1997) Physical status and expression of HPV genes in cervical cancers. *Gynecologic oncology* **65**: 121-129
- Park TW, Fujiwara H, Wright TC (1995) Molecular biology of cervical cancer and its precursors. *Cancer* **76**: 1902-1913
- Peitsaro P, Johansson B, Syrjänen S (2002) Integrated human papillomavirus type 16 is frequently found in cervical cancer precursors as demonstrated by a novel quantitative real-time PCR technique. *Journal of clinical microbiology* **40**: 886-891
- Peter M, Rosty C, Couturier J, Radvanyi F, Teshima H, Sastre-Garau X (2006) MYC activation associated with the integration of HPV DNA at the MYC locus in genital tumors. *Oncogene* **25**: 5985-5993
- Pfreundschuh M, Genth B, Kirchner H, Scheurich P, Lindemann A, Steinmetz HT, Schaadt M, Diehl V (1989) [Mechanism of resistance to tumor necrosis factor]. *Onkologie* **12**: 128-130
- Ramsay RG, Gonda TJ (2008) MYB function in normal and cancer cells. *Nature reviews cancer* **8**: 523-534
- Reuter S (1995) Identifizierung eines neuen Zinkfingergens im Integrationsbereich von Papillomavirus-DNA im Genom der Zervixkarzinomzelllinie ME180 und Charakterisierung der viralen Integratstruktur. Dr. rer. nat. Thesis, Naturwissenschaftlich-Mathematischen Gesamtfakultät, Ruprecht-Karls-Universität Heidelberg, Heidelberg
- Reuter S, Bartelmann M, Vogt M, Geisen C, Napierski I, Kahn T, Delius H, Lichter P, Weitz S, Korn B, Schwarz E (1998) APM-1, a novel human gene, identified by aberrant co-transcription with papillomavirus oncogenes in a cervical carcinoma cell line, encodes a BTB/POZ-zinc finger protein with growth inhibitory activity. *The EMBO journal* **17**: 215-222
- Reuter S, Delius H, Kahn T, Hofmann B, zur Hausen H, Schwarz E (1991) Characterization of a novel human papillomavirus DNA in the cervical carcinoma cell line ME180. *Journal of virology* **65**: 5564-5568
- Robertson G, Bilenky M, Lin K, He A, Yuen W, Dagpinar M, Varhol R, Teague K, Griffith OL, Zhang X, Pan Y, Hassel M, Sleumer MC, Pan W, Pleasance ED, Chuang M, Hao H, Li YY, Robertson N, Fjdl C, Li B, Montgomery SB, Astakhova T, Zhou J, Sander J, Siddiqui AS, Jones SJ (2006) cisRED: a database system for genome-scale computational discovery of regulatory elements. *Nucleic acids research* **34**: D68-73
- Romanczuk H, Howley PM (1992) Disruption of either the E1 or the E2 regulatory gene of human papillomavirus type 16 increases viral immortalization capacity. *Proceedings of the National Academy of Sciences of the United States of America* **89**: 3159-3163
- Saunier M, Monnier-Benoit S, Mauny F, Dalstein V, Briolat J, Riethmuller D, Kantelip B, Schwarz E, Mougin C, Pretet JL (2008) Analysis of human papillomavirus type 16 (HPV16) DNA load and physical state for identification of HPV16-infected women with high-grade lesions or cervical carcinoma. *Journal of clinical microbiology* **46**: 3678-3685
- Schiffman M, Castle PE, Jeronimo J, Rodriguez AC, Wacholder S (2007) Human papillomavirus and cervical cancer. *Lancet* **370**: 890-907

- Schiffman M, Rodriguez AC, Chen Z, Wacholder S, Herrero R, Hildesheim A, Desalle R, Befano B, Yu K, Safaeian M, Sherman ME, Morales J, Guillen D, Alfaro M, Hutchinson M, Solomon D, Castle PE, Burk RD (2010) A population-based prospective study of carcinogenic human papillomavirus variant lineages, viral persistence, and cervical neoplasia. *Cancer research* **70**: 3159-3169
- Schmitt M, Bravo IG, Snijders PJ, Gissmann L, Pawlita M, Waterboer T (2006) Bead-based multiplex genotyping of human papillomaviruses. *Journal of clinical microbiology* **44**: 504-512
- Schneider-Gadicke A, Schwarz E (1986) Different human cervical carcinoma cell lines show similar transcription patterns of human papillomavirus type 18 early genes. *The EMBO journal* **5**: 2285-2292
- Schwarz E, Freese UK, Gissmann L, Mayer W, Roggenbuck B, Stremlau A, zur Hausen H (1985) Structure and transcription of human papillomavirus sequences in cervical carcinoma cells. *Nature* **314**: 111-114
- Seedorf K, Krammer G, Durst M, Suhai S, Rowekamp WG (1985) Human papillomavirus type 16 DNA sequence. *Virology* **145**: 181-185
- Shendure J, Ji H (2008) Next-generation DNA sequencing. *Nature biotechnology* **26**: 1135-1145
- Smits HL, Cornelissen MT, Jebbink MF, van den Tweel JG, Struyk AP, Briet M, ter Schegget J (1991) Human papillomavirus type 16 transcripts expressed from viral-cellular junctions and full-length viral copies in CaSki cells and in a cervical carcinoma. *Virology* **182**: 870-873
- Smotkin D, Wettstein FO (1986) Transcription of human papillomavirus type 16 early genes in a cervical cancer and a cancer-derived cell line and identification of the E7 protein. *Proceedings of the National Academy of Sciences of the United States of America* **83**: 4680-4684
- Snijders PJ, Steenbergen RD, Heideman DA, Meijer CJ (2006) HPV-mediated cervical carcinogenesis: concepts and clinical implications. *The Journal of pathology* **208**: 152-164
- Solomon D, Davey D, Kurman R, Moriarty A, O'Connor D, Prey M, Raab S, Sherman M, Wilbur D, Wright T, Jr., Young N (2002) The 2001 Bethesda System: terminology for reporting results of cervical cytology. *Jama* **287**: 2114-2119
- Spalholz BA, Yang YC, Howley PM (1985) Transactivation of a bovine papilloma virus transcriptional regulatory element by the E2 gene product. *Cell* **42**: 183-191
- Stanley MA, Browne HM, Appleby M, Minson AC (1989) Properties of a non-tumorigenic human cervical keratinocyte cell line. *International journal of cancer* **43**: 672-676
- Steinmeyer N (2009) Mutationen und Integration von DNA humaner Hochrisiko-Papillomaviren in anogenitalen Krebsvorstufen und Karzinomen. Diplom Thesis, Fachbereich Chemie- und Biotechnologie im Studiengang Biotechnologie, Hochschule Darmstadt, Darmstadt
- Stenlund A (2003) Initiation of DNA replication: lessons from viral initiator proteins. *Nat Rev Mol Cell Biol* **4**: 777-785
- Strauss WM (2001) Preparation of genomic DNA from mammalian tissue. *Curr Protoc Mol Biol* **00**:2.2.1-3
- Swan DC, Rajeevan M, Tortolero-Luna G, Follen M, Tucker RA, Unger ER (2005) Human papillomavirus type 16 E2 and E6/E7 variants. *Gynecologic oncology* **96**: 695-700
- Sykes JA, Whitescarver J, Jernstrom P, Nolan JF, Byatt P (1970) Some properties of a new epithelial cell line of human origin. *Journal of the National Cancer Institute* **45**: 107-122
- Thierry F (2009) Transcriptional regulation of the papillomavirus oncogenes by cellular and viral transcription factors in cervical carcinoma. *Virology* **384**: 375-379
- Thierry F, Yaniv M (1987) The BPV1-E2 trans-acting protein can be either an activator or a repressor of the HPV18 regulatory region. *The EMBO journal* **6**: 3391-3397

Thorland EC, Myers SL, Gostout BS, Smith DI (2003) Common fragile sites are preferential targets for HPV16 integrations in cervical tumors. *Oncogene* **22**: 1225-1237

Trottier H, Franco EL (2006) The epidemiology of genital human papillomavirus infection. *Vaccine* **24 Suppl 1**: S1-15

Valente EM, Brancati F, Silhavy JL, Castori M, Marsh SE, Barrano G, Bertini E, Boltshauser E, Zaki MS, Abdel-Aleem A, Abdel-Salam GM, Bellacchio E, Battini R, Cruse RP, Dobyns WB, Krishnamoorthy KS, Lagier-Tourenne C, Magee A, Pascual-Castroviejo I, Salpietro CD, Sarco D, Dallapiccola B, Gleeson JG (2006) AHI1 gene mutations cause specific forms of Joubert syndrome-related disorders. *Annals of neurology* **59**: 527-534

van Es JH, Kirkpatrick C, van de Wetering M, Molenaar M, Miles A, Kuipers J, Destree O, Peifer M, Clevers H (1999) Identification of APC2, a homologue of the adenomatous polyposis coli tumour suppressor. *Curr Biol* **9**: 105-108

Van Tine BA, Kappes JC, Banerjee NS, Knops J, Lai L, Steenbergen RD, Meijer CL, Snijders PJ, Chatis P, Broker TR, Moen PT, Jr., Chow LT (2004) Clonal selection for transcriptionally active viral oncogenes during progression to cancer. *Journal of virology* **78**: 11172-11186

Veress G, Szarka K, Dong XP, Gergely L, Pfister H (1999) Functional significance of sequence variation in the E2 gene and the long control region of human papillomavirus type 16. *The Journal of general virology* **80 (Pt 4)**: 1035-1043

Vinokurova S, Wentzensen N, Kraus I, Klaes R, Driesch C, Melsheimer P, Kisseljov F, Durst M, Schneider A, von Knebel Doeberitz M (2008) Type-dependent integration frequency of human papillomavirus genomes in cervical lesions. *Cancer research* **68**: 307-313

Vizcaino AP, Moreno V, Bosch FX, Munoz N, Barros-Dios XM, Parkin DM (1998) International trends in the incidence of cervical cancer: I. Adenocarcinoma and adenosquamous cell carcinomas. *International journal of cancer* **75**: 536-545

von Knebel Doeberitz M, Bauknecht T, Bartsch D, zur Hausen H (1991) Influence of chromosomal integration on glucocorticoid-regulated transcription of growth-stimulating papillomavirus genes E6 and E7 in cervical carcinoma cells. *Proceedings of the National Academy of Sciences of the United States of America* **88**: 1411-1415

Walboomers JM, Jacobs MV, Manos MM, Bosch FX, Kummer JA, Shah KV, Snijders PJ, Peto J, Meijer CJ, Munoz N (1999) Human papillomavirus is a necessary cause of invasive cervical cancer worldwide. *The Journal of pathology* **189**: 12-19

Wang Q, Griffin H, Southern S, Jackson D, Martin A, McIntosh P, Davy C, Masterson PJ, Walker PA, Laskey P, Omary MB, Doorbar J (2004) Functional analysis of the human papillomavirus type 16 E1=E4 protein provides a mechanism for in vivo and in vitro keratin filament reorganization. *Journal of virology* **78**: 821-833

Watts KJ, Thompson CH, Cossart YE, Rose BR (2002) Sequence variation and physical state of human papillomavirus type 16 cervical cancer isolates from Australia and New Caledonia. *International journal of cancer* **97**: 868-874

Wellmann A, Fogt F, Hollerbach S, Hahne J, Koenig-Hoffmann K, Smeets D, Brinkmann U (2010) Polymorphisms of the apoptosis-associated gene DP1L1 (deleted in polyposis 1-like 1) in colon cancer and inflammatory bowel disease. *Journal of cancer research and clinical oncology* **136**: 795-802

Wentzensen N, Ridder R, Klaes R, Vinokurova S, Schaefer U, Doeberitz MK (2002) Characterization of viral-cellular fusion transcripts in a large series of HPV16 and 18 positive anogenital lesions. *Oncogene* **21**: 419-426

- Wentzensen N, Vinokurova S, von Knebel Doeberitz M (2004) Systematic review of genomic integration sites of human papillomavirus genomes in epithelial dysplasia and invasive cancer of the female lower genital tract. *Cancer research* **64**: 3878-3884
- Werness BA, Levine AJ, Howley PM (1990) Association of human papillomavirus types 16 and 18 E6 proteins with p53. *Science* **248**: 76-79
- WHO/ICO Information Centre on HPV and Cervical Cancer (HPV Information Centre). (2010) Human Papillomavirus and Related Cancers in World. Summary Report Update 2010.
- Woodman CB, Collins SI, Young LS (2007) The natural history of cervical HPV infection: unresolved issues. *Nature reviews* **7**: 11-22
- Wu X, Zhang C, Feng S, Liu C, Li Y, Yang Y, Gao J, Li H, Meng S, Li L, Zhang Y, Hu X, Wu X, Lin L, Li X, Wang Y (2009) Detection of HPV types and neutralizing antibodies in Gansu province, China. *Journal of medical virology* **81**: 693-702
- Xu B (2010) Integration of Human Papillomavirus Type 16 DNA in Cervical Carcinogenesis: Design of a novel strategy for HPV16 integration site determination in cervical scrapes and analysis of HPV16-induced c-myc insertional mutagenesis. Dr. rer. nat. Thesis, Fakultät für Biowissenschaften, Ruprecht-Karls-Universität Heidelberg, Heidelberg
- Yamada T, Manos MM, Peto J, Greer CE, Munoz N, Bosch FX, Wheeler CM (1997) Human papillomavirus type 16 sequence variation in cervical cancers: a worldwide perspective. *Journal of virology* **71**: 2463-2472
- Yamada T, Wheeler CM, Halpern AL, Stewart AC, Hildesheim A, Jenison SA (1995) Human papillomavirus type 16 variant lineages in United States populations characterized by nucleotide sequence analysis of the E6, L2, and L1 coding segments. *Journal of virology* **69**: 7743-7753
- Yu T, Ferber MJ, Cheung TH, Chung TK, Wong YF, Smith DI (2005) The role of viral integration in the development of cervical cancer. *Cancer genetics and cytogenetics* **158**: 27-34
- Zehbe I, Wilander E, Delius H, Tommasino M (1998) Human papillomavirus 16 E6 variants are more prevalent in invasive cervical carcinoma than the prototype. *Cancer research* **58**: 829-833
- Zhang L, Murray F, Zahno A, Kanter JR, Chou D, Suda R, Fenlon M, Rassenti L, Cottam H, Kipps TJ, Insel PA (2008) Cyclic nucleotide phosphodiesterase profiling reveals increased expression of phosphodiesterase 7B in chronic lymphocytic leukemia. *Proceedings of the National Academy of Sciences of the United States of America* **105**: 19532-19537
- Zheng ZM, Baker CC (2006) Papillomavirus genome structure, expression, and post-transcriptional regulation. *Front Biosci* **11**: 2286-2302
- Ziegert C, Wentzensen N, Vinokurova S, Kissel'jov F, Eienkel J, Hoeckel M, von Knebel Doeberitz M (2003) A comprehensive analysis of HPV integration loci in anogenital lesions combining transcript and genome-based amplification techniques. *Oncogene* **22**: 3977-3984
- zur Hausen H (1999) immortalization of human cells and their malignant conversion by high risk human papillomavirus genotypes. *Seminars in cancer biology* **9**: 405-411
- zur Hausen H (2000) Papillomaviruses causing cancer: evasion from host-cell control in early events in carcinogenesis. *Journal of the National Cancer Institute* **92**: 690-698
- zur Hausen H (2002) Papillomaviruses and cancer: from basic studies to clinical application. *Nature reviews cancer* **2**: 342-350

Appendix

A1. Nucleotide sequences of integrated HPV68b in ME180 and ME180R

HPV68b(int)-ME180

The sequence is described in Results section 2.1.1, including Figures 2.3-2.4, and Table 2.1. Black color represents cellular sequence of chromosome 18. Red color and blue color represent the 5' copy and 3' copy of the integrated HPV68b respectively. Underlined nucleotides indicate overlapping areas between viral and cellular DNA.

1 ACACACCAAT GGAGAGAGT GCCACGCGAG ATGGCTGGCA CCTTGGCCGG GTTGGCCATG
61 GAAGTCTGGC GCAGGCAAAAT CAAATGGTGC CCCTTCAAGC TGTGTGTTCTT TAAATATTTT
121 TTGATCCCTTT ATTTCVTRAG GGCAGTGGG GGGAGACTGG TTTTCTTTTA AGAAGAATGG
181 CACCAATGT TTTTAAAC ATGCTTTGG CTTTGAAGC CATAGTTTCA TACTTGTGTC
241 CCTTAATGTC CACTGCTCTG CATATATGC AGTTTGAAGA ATGCTCTCTA GGCACCTGTC
301 CAGTTCAATG GCTTGATATG TTTAGTATTT TGAATGTTTT ACAGGAGGCA CCGCCTGGTG
361 GTCAATATTA TGCATACCAC GTTCTCGTGC TGCATAATAT ATTGCATTTT CCAGTCCGAT
421 ACAGTTTCAA TAGTAAATAT GGTCTTTTAT ACATTTAATG TCCCTGTTCAT AATGTTCTAA
481 TCTTTTCTGC TGTAAACATAT TTAACGTTCG GGAAAGTGTG TTAATCAATG TCTCCCTTCA
541 CTTCTCTCTG CTGCAAGTCT AATTGCGACC AAGTCTTTTC AAAAACAAT TTCCAGATTTT
601 TATCATATGAT TGTATACACT GGGTTCCTGT TTTGCTCAAA TGGAAATGCA TTAGGAATTT
661 TAGACACAGT TAGTCTGCTA TGTAAATAG GCCACCTATG GTCTTCCACGA GGGTGTGTAT
721 TTGATGTTAT TAGCATGTGT GGAATTTTA TTTGATATAT GTGCTGTGCT TTTCTAATCTA
781 AACTATTATGG GTTACACTAT ATGGCTATTT TCGATGTAAT ATCTAAATAT GACCAACATG
841 TACTGTTTGG GTCATCTTAAC ATGGCTATTT TTGCATCTGC AAGTGCTCTT AACCAAAATG
901 GACTAGCTGA ATTTACATAT GAAATATATT TGCCCTTGCAA GAAATGATTA AGACTCATGC
961 AAAAATATGA CTTGCCCTGA TTTGGTGGCC CATGTATAAC TATACAAATTA CGTTTGGCG
1021 TGCTTTTAA AAAAACTCTT AATGCACATA AAAATGTTAT AAATCTTAG CCCTGATATC
1081 TTAAATCATG TACATATGTT CGCCATCATC CGCCTTCATC ACATTTACTG CATCTAATTT
1141 TAAATCATG GGGCATGAG ATTTGTGCTT TTTGTGCGCG TTTGTGAATG CTACACATFG
1201 TTGCACATC TTTTACATAT TTTGCTTGAC AGTTGCTTTT TAAAAACGCT GCAGCATATC
1261 TATTCAATC TGCCACATA GCATATGAA ATGCTATATC ACTTCATCT GTTAACATCT
1321 TATCAATATG CCAATGPAAC ATGCTGATA GATCAAAATC ACTATCACT ATTCCATGTT
1381 GTATATATAT TAAATTTTAT ATCCATCTG GCGTGTGTTT ACACACTCTA CTAAATATAG
1441 ATATTCTCTG TCTATACCAA TACAATGCTG CAAAGGGGCT ACGCAATTTT GGTGGCTGCA
1501 AAAGCATACA GCTGCTGGA ACATGCAACA ATGTACTTAA TCCTTTCTCT ACTGTATCTA
1561 TATTTTCTCC ACATTTGTAT CTTATTAAAC TTAATATTAA TATTCGTTT TTTGTATCTA
1621 AACTTTGTAT ATGGGTATAT AATGCATATT GTTTAATTAG TGTTTTAAAC CCTTCGGCAA
1681 TGGTGTGATT PACTCGGAAT ATTCGTGCTA CCGAGTCCGT GCATGGTGC TTATCACTTT
1741 TAAATGTAG TACTTAGATA TTAAGGGCA ATCCATATAG TTTTGTAAAT TCTGTAAACA
1801 TTGAGCTTT TTTATATATA CATGTAATA ATACTTTTAG TTGCGTAGTA GGTGATTAG
1861 GATCTGTTT TTTACTATAT ACAGACTGT CTACACTACT ACAGTCTCTC CGTATGCTGT
1921 GACCTGTTTC TCCCGGACG CCCCCTGCG CCCCATTGTT ATTAGTTGCT ACAGTTACTC
1981 CCGAGTATAG TTTGCATTTCC ATATTGTATG TTTCTGTAAAT GGGGACTTTG CTTAAAGGGCT
2041 CTTTCTTATA CTCTCTGCTA ACTTTCGTTT TAGGGCAAGC TCTGCTGTTA TACAATATC
2101 GGCCGTGTGC ATATTAAACA GTACGTGCGC TGTGTCAGCG TCTGCTGTTA TCAATATC
2161 TGTAGCATCA TCAATATATC TACATATGTC TGAACATGTA TCTGTGCGGT TTAATATCCT
2221 ATCTCTTAG ACTGTGACG CTGTGTTGTT ATCTACTATT CTTGTACATA AATACATACC
2281 GTTACCTGCC GTCCGPTCC CATCTGTAC TTCAATATG GCAATGACGA ATTACTGGGT
2341 TTTGCTGCTA TCCCGGACG CTTCTACTAC TAGTGTGCAAT AGGTGTTTAC ACTTCAACA
2401 CCGAGTGTG ATGTGTGAC GTTGTGTGTC GTCCCGTCTG GCTGTAGTGT GATGTGCGA
2461 GTGATCACTG CAGTGTGTCG GTTCACTTAT TTAGTCTGCT GAATCTCTCA ATTGCTGCTG
2521 GTGAATAAAG TCGAGTTCGG TCAATTTCAAT CCGATGGACAT AACCTTAACA CAATTTCTG
2581 ACATACAAG TCGAGTTCGG TCAATTTCAAT TTTTACTGTT TTTTACTGTT TTTTCTCTG
2641 CAGGTGGGCT TTTTGGTCCAT GCATGTATG GGGGACATG TCTGTAAAG TTTTCTCTCA
2701 TCGGTCTCTC TCGTTTCACT GTCCAGCATG GCGGCTGCTC TCTGTAAAG TTTTCTCTCA
2761 TTTTATGAAA TCTTCTGTTT GAAATTTAGGT GCTTAGTATT TTACGACGA CTCAATGTTT

2821 TCAGGCAACA CATGACCTT ATTGATAAT CATATAACT CATATACTT TGTATTAGTT ATGTTTCTTA
2881 ATGTTGTTGC ATACACCGAT TCTGAGTAAT ATCGTAGTTC CCGCATTTTC GCTAAATAAT
2941 TAATATATGA TTGGCATGA GCTAATGTTA CCCCGTCCCT ATATACTTCA TTTAAGTTCAC
3001 CAAAGGQAA TTCTATAACC TTGTCCTGTT TGTGTGCTT TCTGCAATAG ACACAGTCTA
3061 TTTCTAACCT GTGCCAATG GTGCCAATG GTGCCAATG GTGCCAATG GTGCCAATG GTGCCAATG
3121 GTTCTCTCAG GTTGGAAT AGGCCAATG GTATAGACA CTGCTGGT CAGCTTTATA
3181 TACACCTGTT TCGGTGATG TACTTCTCT TACTTCTCT TATCTATTA AATGAATTTA
3241 AACAATGTTA AGTAAAGTA TGGTATGTC TCCCAACCT TATTCATGTT CACAATTTTA
3301 TGGATGTGGT ATTAGTCACT GTATACATTA GCTAATCTTG GCATCTCTG ATATGCCATA
3361 TAGTTATACA GGCACACTAT GTATTAGGC AAGTAAACAG AGTTAAACAG CTAATGCCATA
3421 AAGCAGTTTT ATTCAAAGG TGAAGATATA TAGCTGCCA AACTATTGT GCACTGCACC
3481 TGGACAGGAT GATGACTAAG TAGGTGCGC CAGTACTAC TGTTACCGGT GCCACTATG
3541 TGGGTAAACC AAGGTGAAC ATGTTTTGC TGAAGACAT AGTTTTAAT AGTTTTAAT
3601 AGGAATAGC ACCACGACC GAAACGHT GCACAAAA TGGCCITACA AATGGCCAC
3661 TTGTTAAT ATAGAACTAT ATATATTAC TACTCTGTA CACACTTAG GGTAGGGTA
3721 CAAATGAT AGGAAATGTT GAGTCCCTAT GATATGTA ACAAAGCAT ATGTATGTA
3781 CAGTTTGT ACATAGATA CACTAGATA TACTTTAT TAACAAAA AGATATACA
3841 ACATATCCCG CAAATACATA CATATACATA TACATACATA TGCACACAT ACCAACAA
3901 ACATACACG TATACACAC CACATACAC CACATACAC AACATACAA ACACAAATTA
3961 CTTTGACGA CTTTACGTT TGCTGTAGA GTGTAGTGA CTTCACACA GTAAAGATA ATTTGATG
4021 TTTAGGGGG CCTTACGTC GTGTCGCGC GACCTGCC TGTAAAGAA ATTTGCTCC
4081 TAAAGGAAC TGGTCCAGTT CAGAACTAAA CTTTCTCTT AATTTACAT TCCAAAAGTT
4141 TAGCCATCA TATGATCTT TTTTAGT TTTTAGT TGCAGGCG TCTTTTTHG ATGTATTTG
4201 TGTGATTTGC AGATAGCGT ATGTATCTAC AAGACTAGCA GATGTTGGAG GGGCAACACC
4261 AAAAAATCCAA TCACCCAAA TAGCAGGAT CATAGTATG ATATAGACA TTAATACGT
4321 GGAATGTT ATAGTACACA ACTGAAAT AATTTGCAAA TATATACTT CAACATGCTT
4381 AATATATTCC TTAATTTTAT TAGCATGATA AATAATTGTT ACAGTGAAT CAGTAGTAGT
4441 AGATAAATGA AATTTGTTAC TGCCTGTCG CTTGTGCGC CAATAGCAT AGCTACTAGG
4501 CCAACAAATA CCATTTGTT TGCCCTGTC CTTGTGCGC CAATAGCAT AGCTACTAGG
4561 TAGGAGTCT GAGGATACA TAGCCACT AGGCGAGGG CAAATACAT AACTACTAGG
4621 ACTGTACGT ATGTACGTC CCTTATATA CAAATCAGTA GGTATAGT CCGTACCAT
4681 GCCCTCTCA TTCAAAAAT CAGTACAAA TAACTAGTCC GATGTAAAC AAGAAACAT
4741 ACTGCTCTCA TATACATCTG CACATATCTG CCAATAGTCA GATATTTGC AGACTGATG
4801 ATATATCT TAAAGCACT CCGTTTTGT TCTTGTAAAT GTTACTAAAT CCAATAGCT
4861 ATATCTGTA TCAATCATAT CGCATCTGT AATAGTGTG TTTACTAAT CCAATGTGG
4921 AAGTTCGCG CGTGCACAT TGTAGGCTT ACAAGATTTA CTTTGGCCC AGTGTCCCC
4981 AATGGCAGC ACACAGCCTA TAATACATG TTTGCTTTGT TTAATGTTA CTAACAACT
5041 GCTCTACTG TCTTTAGGAT TTTTGTGGA GGAACCGGG GAATTTCCG TATCATCTAG
5101 CACTATATAT ATGGAATGCC CACTAAGGCC AACTACTAAT GCGTCCGCC TACTATTTTC
5161 AACACAAACA CAGGCCATA CCAATGCTG CTGATCAGG TTATATAG TAGACTCAGG
5221 AAGCTTAAT TTTATAGAT CAGTATGGA AHTCTAAC ACCCTGAT TATATAGTA
5281 AAGCTTAGA ATGCTCTGCT TGCGGCCCC AGACATAGG ACCTTAAAT ATGATGGCC
5341 TACAGTATG ACCACTATG TACCAGATA TACCAGATA CCAATAGAT CCAGTCCGT TTAAGTAACT
5401 ATCTGTATG AACACTGTC CACTCAGGG GGGAGCAAA TACACATGT TGTGCTAGA
5461 GCSCACAAAT GCCATGCA AAAAAAAG GAAGCGTGT ACCTTTTTT AATAAAGA
5521 ACAATAATG TAAATAATA TAAATGTTGC CATATAGT TATGCTAG TTTGATACAA
5581 TGTGTTAGA GGGTGTAAA GGCACTGTC GAGCGTTGC TGTAAACACA ACATCAGC
5641 CAGTATTTAC AGGCTGCTC CAAAGATAG CAAATAGGAT AGTAGTGA GATATGNG
5701 TAGATGCTGT AGAGCTAAT CAGGACAG CAAATAGGAT AGTAGTGA AATGTAGCT
5761 TATGAAATG AGTATCCAT ACTGTAGTAT TGTAGTATC TGTAGTATC ATCATATTA
5821 AAGTATCCAT AGGGTACAG TGTCTGGG CAAACAAAG TTGTAGTCA ATGCTATGA
5881 CAGGAGCAAT GCCATAATA TCAATGAT ATGTCACTG TGCCCAAT TGTGATCCC
5941 GCGGTGTAAA CATATGTC TTTTTCCTA CTCTGTAAA ACGTATGTC CTCTTCGGG
6001 AAGTTAAGCG AGGCTATGT AAACGAACA TGTCCAGAAA ATCCGATCA GGAGTATGT
6061 CAGTAGTTC ATATGTAAT GTAGTATCAA CAGGCTCAA AATTAAC AGGAGATTA TCAATGTTA
6121 CAAATATGA AGGTTGAGTT ACAAATCAA AATTAAC AATTAAC AGGAGTGA CTACTAAG
6181 TACTATATG AGGTTGAGTT GCGACAGG TAAACCCAG TATAGTGA CTACTAAG
6241 GTTCGTACC AGTGCATGT GTGCAATA CTGCAATG CTTCTCTCA TATCAATG
6301 TCTCCGATG GGGGTACT ATAAACAT TACAGAGAC TTAACCTGT TTAGGACTT
6361 CTATATAGT GGGGTCTGCA AATGAGGCT TAGTAGGCT AGTAAACT AGTACTGT ACTTGCACG
6421 ACCAGACAG TCTAACAG CAGGTGATG AGTGTAGT AGTGTAGT AGTGTAGT
6481 CAAACGAGA AGTGTCTGTA TATGTTGTA CCGGTGTC AGATGATA ACCTGAAAT
6541 TGTCCCAA GAGGTTCTG TAGGACCCAC AGGTTCAATA ACCACAGGT
6601 CAGGTGCAAG CCAACATCT ACACAGTAT TAGTGTACC ACCTAAAGG ATGTACCGAG
6661 TACGACCCC GGTCTGTGAC CAGTACAA TGGCTAGGCC ACCCAAAA ATACCTAAC
6721 TGGTCCATG CAAATGTTG TCTGAGAGT TGGTCTCTC BACCTTAAT ATACATCAG
6781 GAGGACATG CCGTATGTT TATATATC AGTTGACAT GCACGTTGC
6841 CCGTGCAGC AGGTGTGAT ACATATTTA TTTACAAAA TATACACCA TACAAATGA
6901 CAATCTACT TATTAGCAT TACATAGA CAAAACCTG GTAAATACCA CACAGGCT
6961 AAAAAAAA GTATATATAC AGCAAGACC TCCAATGGTG TGTAGCTAGC TAATATAAC

11221 TTGCAGAAAG CAACTACAC GCACGAGGT ATATGAATTT GCCTTTGGTG ACTTAAATGT
11281 AGTATATAGG GACGGGTAC CATTAGCTGC ATGCCTAATC TGTATTAAAT TTATCATGAA
11341 AATATGCGAA CTACGATATT ACTCAGANTC GGTGTATGCA ACACATATTAG AATACGTAAC
11401 TATATACAAAG TTATATGATT TATCAATAAG GTGATCTGTG TGCTCTGAAC CATGTAGTCC
11461 TGCTGAAAAA CTAGAGCACC TAAATTCAAA AGCAAGATTT CATAAAATAG CAGAAAGCTT
11521 TGACGACAG TGTGGCATT CTGGGACAG TAAACGAGG GACCGACAG CATAAAATAG CAGACAGGA
11581 GAAACACAA GTATAAACTA ACTATGCATG CAGTAAAGCC CACCGTGCAG GAAATTTGGT
11641 TAGAGTTATG TCCATGCAAT GAAATAGAG CGGTGCACTT TGTATCTCAC GAGCAATTAG
11701 GAGATTCAGA CAGTGAATA TAATACCCG ACCATGCAGT TAATCACAC CAAACATCAAC
11761 TACTAGCCAG ACGGACGAA CAACAGCCTC ACACAATTCA GTGTAGCTGT TGTAAAGTGA
11821 ACAACCTACT CCACTAGTA CTAGAAGCT CCGGGAGAA CCTCGGAAG CTAGAATGCG
11881 TGTTTATGGA CTCACAAAT TTGTCTGTG CFTGTGTGC AACGAAAC CAGTAATCTG
11941 CAAATGCGAA TTGTGAAGGT ACAGATGGG ACGGAGGG GTGTACGGA TGGTFTTTTG
12001 TACAGCAAT AGTAGATAA CAACAGGTG ACAGATGGG ACGAGTCTC AGAGATAG GATGAAAACG
12061 AGCAGATAT AGGTTTCAGC ATGGTAGATT TCAATTGATG TGTACAGTA ATTTGTATAC
12121 AGCAGAGCG TGACACAGA CAGTACTGT TAAATATGCA ACAGGCCAA AGGATGGAC
12181 AACAGTGG TGCCCTTAAA CGAAATGAT CAGACAGTAT AGAAAGCAG CCGTTAGCAA
12241 AGTCGCAAT ACAGAACTA TATGAAATG TAAGCAGTAC ACAGGACAGA CAACCGCGT
12301 ATACAGTCC GGACAGCGG TATGCAATA TGGAAATGGA AACTAATCG GAGGTAACTG
12361 TAGCACTAA TACAAATGGG GCGAGCGGG GAAATGAGG GAAATGCG GACACATAC
12421 GCGAGGACTG TAGTAGTGA CACAGTCTA TAGATAGTGA AAACAGGAT CCTAATCAC
12481 CTACTAGCA ACTAAAAGTA TTATTAATAT TATTAATAT AAAAGCTGCA ATGTTAACAG
12541 AHTTAHAA AGTATATGA TTGTCTTTA ATGACTTGT ACATCAAT AATGAGTAA
12601 AGACACAT TGCGACTCG GTAGACGAA TATTCGAGT AAATCAAC ATTGCCGAAG
12661 GGTCTAAA ACCTAATAAA CAATATGAT TATATACCA TATCAATGT TTAGATACAA
12721 AAAACGAA ATTAATATA ATGTATAA GATACAAAT TGGGAAAA AGAATAACAG
12781 TAGGAAAAG ATTAAGTACA TTGTGCAAT TTCCAGACAG CTGTATGCTT TTGCAGCCAC
12841 CAAATTTGCG TAGCCCTGTT CCGAATGTT ATTTGATAG AACAGATA TCTAATATTA
12901 GTGAGTGTG TGGAGACAG CCAATGAGA TAAAAGATT AACTATAA CAACATGAA
12961 TAGAGTATG TGTATTTGAT CTATCAGCA TGGTACATG GGCATTTGAT AATGAGTTAA
13021 CAGATGAAG TGAATAGCA TTTCATATG CTATGTTGCG AGATTGTAAT AGTATGCTG
13081 CACGTTTTT AAAAGCAAC AGGGGACAA AAACGACAA TGTCAATGCC GCAATGGATT AAATTAGT
13141 GCAGTAAAT TGATGGAGCG GTGATTTGCG GCATGGACT CAGAAATTTAC TACTTGTGTT
13261 TATTACAGG GATTTCTGGA GCCAGA

7021 ACAACACAA GTATCCACAC ATACACACAC ACATGCATGG ACTGCGAAG CCGGACAGTG
7081 CCAATATAT ATATGACAC CCATGACAC ACCAAATAA TGTGTACATG TCCCGTCTCT
7141 TGATGTGTTG CATGATGCC AGTAATGTA TGTGTACATG TCCCGTCTCT
7201 GCAGTCCAG CATGCTGTTG TGTATGTTG GATCTGTTTA TATATATAG
7261 TACACACGA TTGAGCTGTT TGTGTAAT TATATTTTTT TTTTACTGCG CTAGTGGGT
7321 ATACACAGT TTGAGCTGTT TATGATGCC TTAGTCTGTT TATGTATGAT TGTATTTGTT
7381 GTATATTTTT ATAAATATAT ATGATATCAT ACCGCTGCTG CAGGCCAGT GGTGATCTG
7441 CAACTGAAT ATATAAACA TGCACATCAT CAGGACATCG TCCCTCTCAT GTTATAAAT
7501 AGTTGTGAGG CACCACACT CTGACACAA TATTTGCAATG GACCACTTTTA GTTATTTTTT
7561 TGGTGTGCGT AGGATGCTG ACTGGCTCAG GAACGGGGT GTGTACTGGG TACATCTCTT
7621 TAGTGTGTA ACCTAATPACT GTTGTAGAT TTTCGCTGCG ACCTCCATCT GTGTATTATG
7681 AACCTGTTG TCCATGATG CCCCATG TTGCAATGGT GGAAGATCC AGNTATTAT
7741 CACTGTCAG ACGGTFACCA ACATTTACAG GCACCTCTCG GTTTGAAAT TTACATCTTT
7801 CTACACATC ACTTCTGTT TTAGACATTA CCCCCTGCT TGGGCTGTG CAGTATAGC
7861 GATAGTATTT GTTGTCCCA GACGATCTG ACCCTATGGA TCTTCTGAC AGATGAAAC
7921 AATACCTCT TCGAATGTT TTGTAAATG CACCAATPACT ATACTGCAT CCAATAGCT ACATTTACT
7981 GAAACACTT GCATATGTT TGTACACATG GCACTGGTAC AGAACCTAT ATAGTACAC
8041 CTACTCTGTT GTTATGCTGT CTTCTGTTAG CACTGTATG TAGTAGGGCA CATCAACAG
8101 TTCTGCTTAG TATTTGTTAT TTGTAACTC ACCCTTCATC ATTGTGAACA TTGTAAATC
8161 CTGCTTTTGA GCGTGTAT ACTACATTA CATATGAAC TGTGACATA GCTCTGATC
8221 CCGATTTTCT GGCATTTGTT CHTTACATA GGCCTGCTT ACCTCCCGA AGAGCACAG
8281 TAGCTTTTAC CAGATATGCG AAAAGGCAA CTATGTTTAC ACGCGGGGT ACACAATTG
8341 GGCACAGGT GCATATAT CATGATATTA TGTGATGTT TCTTCTGAC AGATGAAAC
8401 TACACCTTT GTTGTCCCA GACGATCTG ACCCTATGGA TCTTCTGAC AGATGAAAC
8461 CACAGATAT TGACATPACT ACAGTATGAG CACTGTATG TCAATAGCT ACATTTACT
8521 CCGTTCCTCA TATATCTGTT CTTCTGTTAG CACTGTATG TCAATAGCT ACATTTACT
8581 CTACTATPCC TATTTGTTAG CTTTGAACA CCGCTGTAAA TACTGTCTCT GATTTGTTG
8641 TACAGCAAC GTCTCCACAG TTGCTTTTAA CACCTCTCAC AACATGAT ACACCTGAT
8701 CCAATACAT ATATGTCAC TATTTATTT TATTAACCAT ATTTGTTCTT TATTAATAA
8761 AAGTAAACG CTTTCTPACT TCTTTGTCAG ATGCAATGTT GCGCTCTGAG CAGACAAATG
8821 GTGTATTTGC CTTCCGCTCT AGTGGGAG AGTTGCAATA CAGATGATA CTTACACCG
8881 ACTGCAATC TGTATGTTG TGTGATCTC AGGTACTTCA CTGATGCTCA TGTAGGCTCA TCAATTTT
8941 AGTCCCTTA TTTCTGGGG CCGCAGCATG GACATCTTCA AGTGTGCGA ATACAATTT
9001 AGGTTGTTA GATTTCTCT ACCTATCTT AATAATTTA GTTCTCTGTA GCTACATTA
9061 TATPAACCTG ATACGACCG ATGTGTATG TGTGTTTGA ATTTGAAAT AGGTAGGGG
9121 CAGCATATG GTTGTGTTAG TGTGGGATC CACTATATA ATAGGCTAGA TGAATCTGAA
9181 AATTTCCCGT TTCTCTCCA CAATAATCT AAGACAGTA GGGCAATGT TTAGTGGA
9241 TATAACAAA GCAACATG TATATAGCG TATATAGCG CCAATGCGA GCATGAGCG
9301 AAGATTAAT CTTCTGATG TACTATGCG CAGCCGGG ACTGTCCAC ATTTGAAATTA
9361 GTAAATAGC CTAATCAGA TGGCGATG ATTGATACAG GATATGTCG TATGACITTT
9421 AGCATATAC AHAACAAA AGCGAGGT CTTTATGATA TATGTCAATC AGTCTGAAA
9481 TATCTGACT ATTTACAAAT GTCTGAGAT GTATATGAG ACAGTATGT TTTTGTTTA
9541 CPTAGGAAC AGTTATGTT TAGGCAATTT TGAATATAG GGGGATGTT AGGGACACT
9601 ATACCTATG AATTTATAT TAAAGTATG GACATACGT ACAGTCTAG TATGTATGTA
9661 TATGCCCTCT CCGTATGTTA TCTTATGTTA TCTTATGAT CCGAGTTAT TAACAAGCC
9721 TATTTGCTG CAAAGGACA GGGACACAC AATGTATTT GTTGGCATTA TCAATTTT
9781 CTTACTGTTG TGGATACAC TCGCATAC AATTTACTT TGTCTACTAC TACTGATCA
9841 GGTGACAA ATATTTAGCA TCCATATAA TTTAAGGAT ATATAGGCA TGTGTGAAA
9901 TATGATTTGC AATTTATTT TCAATGTTGT ACTATACAT TGTCCACTGA TATATGTC
9961 TATATACATA CTAATGATCC TGTATTTTGT GATGATGGA ATTTGGTGT TGGCCCTCA
10021 CCACTGTTA GTCTTGTAGA TACATACCG TATCTGCAAT CAGCAGCAT TACATGCAA
10081 AAGAGCCCG CTGCACATC TAAAAGAT CCAATATGCT GCTTAAACTT TTGGAATGTA
10141 AATTTAAGG AAAAGTTTAT TCTGAACTG GACCAGTTTC CTTTAGAGCG CAAATTTCTT
10201 TTGCAGCAG GTGTCCCGC ATAGCCCACT ATAGCCCGC GTAAAGCCG TGCCACAGCA
10261 ACTACTGAT CTACTCTAA CACAAAGCT AAACGTGTT CAAAGTAAAT GTTGTATGTT
10321 TGTGTTGTTA TGTGTTGTTG ATGTGTTGT GTTATATGTT CATGTGTTGT TGTGTTGTT
10381 GTGCATGAT GTGTATGAT ATATGTTGT GTTTGAGGT ATGTGTTGT AATCTGTTT
10441 TGTATATAA GATATGAT CATTACTT TGTGTTGTA CCGTGACT ACATATGTC
10501 CTTGTTTTAC ATATCATAG ATGCAACAT TTCCATACATA ATTTGAGCG CTACCTAAG
10561 GTGTGTTACA GTACATGTA TATATATA TCTCTAAT ATACCAAGT GCCATTTGT
10621 AAGGCAATTT TGTGTGAC CATTTCGTT CCGTGGTGT ATTTCTCT ATACAGTAT
10681 AAAAATATG TGTGTCAGA AAAACATGTT TCACTTGTGT TTACCCAT AGTTGCAACC
10741 GPTACAGTA TGTACTGCG CACTTACTT AGTCATCAT CTGTCCAGT GCAGTCAAC
10801 AATAGTTTGG CAGCTATAT ATCTCACCC TTGTAAATA ACTGCTTTTA GGTATAGTT
10861 TTTTACTGTT TTTACTGTC TAATAGCAT GTTGCCCTGT TTACTACTT TTGCAATCAA
10921 GATCTGCT TGTATGTC GATATGAT GACTATACC ACATCCATA ATTTGTGCA
10981 CCAATATAG TTGGGCAC ATACCAATAC TTTTACTTAT AACATTTAC AATATTTTA
11041 TAGTATAAG GAGTACCG AAAAGGTC TGAATGAAA CCGTGTATAT AAGCTGAAC
11101 ACAGAGTTG TCTATACAA TGCGCTATT TCAACACTT GAGGACGCG CATCAAAAT
11161 GCAGACCTG TCGAGGACAT ATTCAGTAC GTTACAATAG ACTGTGTCTA

HPV68b(int)-ME180R

The sequence is described in Results section 2.1.2, including Figures 2.6-2.8, and Table 2.3. Black color represents cellular sequence of chromosome 18. Blue color represents the integrated HPV68b DNA. Underlined nucleotides indicate overlapping areas between viral and cellular DNA.

1 GGAATTATGC CTAATTGACA GATGGAAA CTAAAGTGAA GGTGAGGCTA AGGTACAAAA
61 GTGATGGAGT GAGATGGGG GGCCTCGCAT CTTCCTCTGAG GCCATTGGCTT TTGCTTAATGC
121 ACATGTCAGG GGAATTGCGAA GAAAGACAGA AAAAGTGGGT TCCCTACGGG GGTGCTCTGGC
181 CTAAGTCACA AAGTAGGGC TCCTTTCAAGT TTGCTAGGGT GCAACTGACG GGTCAAGAGC
241 TGCGGGGACA CACCAATGCA GAGAGTGCC ACCGCAGATG CTTGCGACCG TGCGGGGGTT
301 GGCCATGGAA GTCGTGGGCA GGCAAATCAA ATGTGTGCCC TTCAAGCTGTT TGTCTTTTAA
361 TTAATTTTTC ACTCTTTAT TTCTAGGGGG CAGTGGGGG AGACTTGTTT TCTTTTAA
421 AAGTGGGAC ACCAATGCTT TTAACAACTT GCTTTGGCTT TGTATGCCAT AGTCTAATTC
481 TTGCTGTCCT TAATGTTCCAC TCTCTCGCAT TATAATGCAAT TTGAGCAATG CTTCTAGTTC
541 CTTAGCTGAG TTCAATGGCT TGATATGCTT TAGTTTTTGA AATGTATTACA GAGGACACCA
601 CCTGTGGTGC ATGATATAGC ATACACGTTT CTCGTGCTGC ATAAATATAT GCATTTTCCA
661 GTCGTATAAC GTTCCAAATG TTAATATGCT CCTTTATACA TTTACTGTCC TGTTCATAAT
721 GTTCTAATAT TTTCTCTGCT AACACATTTA AAGGTTGGGA AAGTGTTCCT ATCAATGTCT
781 CPTCACTCCT GTCCTGTGCT CAAGTCTAAT TTGCAACCAAG TCTTTTCAA AAAACATTC
841 CAGTTTTTAT CATTGATGCT AAGTCTGGG TTCTGTGTTT GGTCAAATGG AAAACATTC
901 GGAAATTTAA ACACAGTTAG TCTGCTATGT AAATACGGCC ACCATTTGTC TTCCACAGTTA
961 TTGCTATTAG ATGATTTTAG CATTGTGGGA CATTTTATTA TGTATTGATG TCTGTGTTTT
1021 CTAATTAAC TTATTTGGTT ACCATCTAAT GCATTTCTTA TGTAAATTAT AAAATATGAC
1081 CACATGTGAC CTGTGTGGTC ATCTACATG GCTATTTTTG CATCTGCAAG TGGCTCTAAC
1141 CAAAGTGAC TAGCTGAATT TACATATGAA ATTAATTGCT CTTGCAAGAA ATGTATAGA
1201 CTTGCGGAAA AATATGACTT GTCCTTAAAT GGTGGCCCAT GTATACTAAT ACAATTGAC
1261 TTTCGCTGCG CTTTATAAAA ACTCTTAAAT GCACATAAA ATGTATAAAT TTCTAGCCCT
1321 TGATATCTTA AAGATGTGAC AATGTGTCG CAATCACCGC CTTCACTACA TTTACTGCAAT
1381 CTAATTTTAA TCCATTTGGG CATTGCAATT TGTGCTTTTT GTGCCGCTTT GTAATGTCTA
1441 CACATGTG CACATCTTTT CATTATTTT TGTGACAGT TGCCTTTTTT TGCCTTTTTT AAGCTGTGCA
1501 GCAATCTAT TACATATGCT CACATAGCA TATGAAGATG CTATATCACT ATATCTAGTT
1561 ACTCATTAAT CAAATGCCA TTGTACCATG TCTGTAGAT CATATCACT CAATCTAATT
1621 CCAATGTTGA TTATGTTTTC TATTTTATC CATTTCTGGG TGTTTTTCACA CACTCTACTA
1681 ATATATAGTA TTTCTGTTCCT ATACCAATAC AATGCTGCAA CAGGCTGCGA CAATTTTGGT
1741 GGCTCGAAAA GCATACATG CTGTGGAACA TGCAACAATG TACTTAAATC TTTCCTACT
1801 GTTATTTCTA TTTTTCACCA TTTGTATCTT ATTAACATTA ATATTAATAT TCGGTTTTTT
1861 GTATCTAAG ATTCATATG GGTATATAAT GCATATGTT TAAATAGTGT TTAAACCTT
1921 TCGGCAATGG TTGATTTTAC TCCGAAATAT TGGTCATTA AAGCAATC CATATCTAT CBTGCTCTTA
2041 GTTAACATTT CAGCTTTTTT ATATTACAT TGTAAATAA CTTTTAGTTG CGTAGTAGTT
2101 GATTTAGAT CTTGGTTTTT ACTATCTATA GCACTGTCTA CACTACTACA GTCCCTCCGT
2161 ATGCTGTCCG CATTTTCCCT TTCACTCTCC CCGTCCGCCC CATTGTGTAIT AGTTGCTACA
2221 GTTACTCTCG AGTTAGTTTC CACTTCCATA TTGATAGTTT CTTGTAATGGC GACTTTGCTA
2281 AAGGGTGCT TTCTATATCT TCTGTATACT TTTCGTTTTAG GGCACGCACT GTTCTGTGAT
2341 CCGTCTGGC CTGTATGATA TTTAACAGTA CCTGTGCTGT CTCACGCTCT GCCTGTATAC
2401 AATATCTGT AGCATATCA ATGAATCTA CCAATGCTGA ACCGTATCT GTGCGTTTTT
2461 CACTCTATC CTTGAGATG GTGTACATG TTGTTTATG TACTATTGCT ATATTTGCT TGTACAAAA
2521 ACTATCTGTT ACACCCGCTC CCGTCCCAT CTGTACCTTC ACATTTGGC ATTCAGATT
2581 ACTGGGTTT CTTTGCACAT CACGACACA CAAATTTAG TGAATCCATA AACACAGTT
2641 CTAGGTTCCG CAGGTTCTCC CGCAGGTTT CTACTACTAG TTGCAGTAGG TTGTTACACT
2701 TACAACAGGT ACATCTGCTT CCGTCTGCTC CCGTCCGCC CATTGTGTAIT AGTTGCTACA
2761 GTTGTGGTGT ATTAATGCA TGTGCGGTTT CATCTATTTT ATGCTGTGAA TCTCTTAATT
2821 GTCGTGACCA TACAAGTGC ACCGCTCTA TTTCAATGCA TGGACATTAAC TCTAACACAA
2881 TTTCCTGCA CAGTGGCTTT GGTCAATGCA TAGTTAGTTT ATACTGTGT TTCTGCGGT
2941 GTGCGTCTG GTCCTCTCG TTTACTGCTC CAGCATGCG GACATGTGC TGTAAATTT
3001 CTTGTTTTT TATGAATCT GCTCAATCT CCAATCTGA ACCGTATCT GTGCGTTTT
3061 AATGTTTTCA GGCACACAT GCACCTTAT CACCGATTCT GAGTAATCAT ATACTTTTG ATAGTTATG
3121 TTTCTAATG TTGTTGATA CACCGATTCT GAGTAATCAT TGTATTCCTG TAGTTCCCG TATTTCCGA
3181 TAAATTTTAA TACATATGTA CACCGATTCT ATGTATACC CGTCCCTATA TACTACATTT
3241 CAGTCACCAA AGGCAATTC ATATACCTCT GTCCGTTTGA GTTGCCTTCT GCAATAGACA
3301 CAGCTATTTG TAACTGTATG CAAATGGTGT TCCATGTCC TGCACAGGTC TGGCAATTTG
3361 TATGCGCGTT CTTCAAGGTT GTCAATGCTA TGCATGTTGA TAGCAACTG CTTGCTTTAC
3421 CTTTATATAC ACCGTTTTCG GTCAAGACCG TTTTGGTGA CTTCCCTTAT ACTATAAAT
3481 GATTTGAAAA TGTATTAAGT AAGATATTTG GTATGTGTC CCAACCTATT TCGGTTGCAC
3541 AATTTATGAG ATGTGGTATT TACTACTGTA TACTATCAC TACAGACAC ATCTTGAT
3601 GCAAAAGTAG TTATAGAGCG CAACTATGCT ATTAGCGAAG TAAAGACAGT TAAACCTTA

3661 TGCTTAAAG CAGTTTTATT ACAAGGGTGG AGATATATAG GCTGCCAAAC TATTGTGCA
3721 CTGCACTGG ACAGATGAT GACTAAGTAA GGTGGCCAG TACATATGT TACCGGTGCC
3781 AACTATGGG GTGAACATG GTGAACATG TTTTGTCTGA BACACATGT TTTTAATCT
3841 GTATAGAGG AATATAGCC ACCGACCGA AAGGTTGCA CACAAATGG CTTTACAAA
3901 TGCCACCTTG AATATATTA TAAATATTA TATATTACAT GTACTGTAC ACACITTAGG
3961 TGAGGCTAC AATATATGA GAAATGTTG CAGTCTATG ATATGTAAA CAGAGCATA
4021 TGTATTACAC AGGTGTCAC CACAAGTAA ACTGACATC ATACTTATT GCACAAAACA
4081 TATTATACAA ACATACCTGC AACATACAT ATATACATC ATACTTATT GCACAAAACA
4141 CCAACACAAA ACATACACAT ATCAACACC AACATACAC ATATACATC ATACTTATT GCACAAAACA
4201 CACACACGG GACAGTCCA CATACATATA CTGTATGAT TATGGATACA TACCACATA
4261 CAGTACATA GTATTGACT GTTATTGTT GTTGTGTTT GTGTGTGAT ACTTATATGT
4321 TGCACTGCC CGCTCTGCA GTCATGATG GTGTGTGTT GTGTGTGAT ACTTGTGTT
4381 TGCTTTTAT TATGATGAT CACAGCTTT GAGGCTTTG CTGTATATAT ACTTTGTTT
4441 TGTGCTCTA TGTAGGTTT ACACGTTT GTCCTTATA GTATGCTATA AGTTTGTTT
4501 TGTGATTTG TATGAGTTA TATTTTATA AATAAATATG GTATCACCC AGTTTGCCAG
4561 TGAGGCTAC GACTCTCAA TAAAATATG TAAAATATG CACAACTAT GCATATGTC
4621 TCCTAGTGT ATAAATAGG TTGAAGGAC CACTCTGCA GACAATATG TCAATGGAC
4681 TGCTGTAGT ATTTTGTGG TGCGCTAGG CATTGGTACT GGTGAGGAA CCGGGTCCG
4741 CACTGGTAC ATTTCTTTAG GTGTATAAC TATACTGTT GTAGATGTT CCGCTGCACG
4801 TCCACGTTG GTTATGAAC CTGAGTTC CAGCAACC TCCATGTGC AATGTGTGA
4861 AGATTCAGT GTTATACAT CTGGCACCC GTGACCAACA TTTACAGCA CTTCTGGGT
4921 TGAATTTACA TCTTCTCTA TCCATACACC TGCCTGTGTA GACATTAACC CTTGCTGCG
4981 GTCTGTGCA GTAGTAGTA CTAGTTTAC TACCTTGA TTTGACACC CACATCTAT
5041 AAGAGTCCCT CAAACAGTG AAGTCTTG TAAATGTT GTAAAGTACC CCAATCGGG
5101 AACACATGGA TATGAGAA TACCTATGA GGTCTCTG TTTATGACT CACTAAATT
5161 TGTGTGTCG TGTGTGCAA CGGAACCCA GATCTCTGA ATGGCAAT GTGAAGTAT
5221 AGATGGGAC GGGAGGGT GTACCGGAT GTTTTGTGA CAAGCATAG TAGATAACA
5281 AACAGTTTC ACAGTCTAG AGATGAGA TGAACCGCG ACAGATACG AGACAGACA
5341 GGTAGTTTC ATTTATGAT GTACAGAT TTGTATACG CGACAGCTG AGACAGACA
5401 GGTACTGTTA ATATGCAAC AGGCCAAG GAGTCAACA ACAGCTGCG CCTAAACG
5461 AAGATTTACA GAGATATAG AAGAGGCC TTTTACAAAG TCGCATTTAC AGGACTATC
5521 AATAATATG AGAGTACAC AGCAAGACA ACCGCGTAT ACAGTCCG ACAGCGGCTA
5581 TGCAATATG GAATGGAAC TAACTCGGA GGTAACTGTA GCAACTATA CAAATGGGC
5641 TGACGGGAG GATGAAGGG AATGGCGA CAGCATACG GAGGACTGTA GTAGTATGA
5701 CAGTGTATA GATGTGAA ACCAGATCC TAAATCCT ACTACCAAC TAAAGTATT
5761 ATTCAGTAT AATATAAA AAGTCAAT AATACAGAA TTTAAAAAG TATATGATT
5821 GTCTTTTAT GACTAGTAC GTACATTTA AGGTGATAAG ACCACATGA CGACTGGT
5881 AGCACAATA TTGACGATA ATCCAACAT TGCCTGAGG TTTAAACAC TAATTAACA
5941 ATATGCAITA TATACCCATA TACAATGTT AGATACAAA AACGAAAT TAAATTAAT
6001 GTTAATAGA TACAAATGTT GGAATAATG AATACAGTA GGAAGAGAT TAAGTACATT
6061 GTTCAATGT CAGACAGCT GTATGCTTT CGACCAACA AATGCTGTA GCTGTGTC
6121 AGCATTTAT TGTATAGA CAGGAATATC TAAATATAG GAGTGTGTG GAGACAGCC
6181 AGAATGATA AAGATTTAA CTATAATCA ACATGGAATA GATGAAGTG ATATTGCT
6241 ATCAGACAT GTACAATGG CATTTGATA TGAATTAACA GATGAAGTG ATATTGCT
6301 TTTATATCT ATTTGGCAG ATTTGGCAG TAACTGCTGA CGTTTATAA AAGCAACTG
6361 TCAAGCAAAA TATGTAAAG ATTTGCAAC AATGTGAGA CATTAACAA GGCACAAA
6421 ACCGAAATG TCAATGCCG AATGANTAA ATTTAGATG AGTAAGTGT ATCAAGCGG
6481 TGAATGGCG ATGACTGCA CAATTTACTA CTTGTGTTTA TTACAGGCGA TTTCTGAGC
6541 CAGA

A2. Nucleotide sequences of HPV68b genomes in CIN2 sample

HPV68b-CIN2

The sequence is described in Results section 2.1.4, including Figure 2.1.4, and Table 2.5:

1 ATGGCGCTAT TTTCACATCC TGAGGACCG CCATACAAAT TGCAGACT GTGCAGGACA
61 TTGGACACCA CATTGCATCA GACTGTGCT ATTGCAGAG GCAACTACAA
121 CGGACAGAGG TATATGAAAT TGCCTTTGGT GACTTAAATG TAGTATATAG GAGCGGGGTA
181 CACTTAGCTG CAGCGCATC ATGTATATAA TTTTATCGCA AAATACGGGA ACTACGAAT
241 TACTCAGAT AGGGTATCG AACACATTA GAACCATATA CTAATGATCA GTTATATGAT
301 TTATCAATTA GGTCATGAT TGCCGTAA CCATAGATC CTGCTGAAA ACTAAGGCAC
361 CTAAATTCAA AACAGATGT TCATAAATA CGAGAAACT TTACAGACA GTGTCGCAC
421 TGCTGGACCA GPAAACAGA CGCACCGCA CGCAACCGG ATGAACACA AGTATAAACT
481 ATGAATGAT GGACAAAGC CCACCGTGCA GGAATATGTG TTCCATGCAA
541 TGAATATGAT CGGCTTGACC TTGTATGCA CCAAGATTTA GGAGATACG ACAGTGAAT
601 AGATCAACC GACCATGAG TTATACCA CCACATCAA CTACTACGA TGCATAGT
661 ACACAGCGT CACACATTC AGTATACGT TTGTAACTGT AACACCTAC TGCACATAGT
721 TTGGACGCG CGCGGAGA ACCTGGGAA GCTAGAATG CTGTTATAG ACTACTAAA
781 TTTTGTGNGT CCGTGGNGT CAAGGAAC CAGTAACTC GCAATGGCA ATTGTAAGG
841 TACAGATGG GACGGACGG GGTGTAAAG ATGGTTTTT GTACAGCAA TAGTATATA
901 ACACAGATG GACGACGCT CAGAGATGA GGATGAANA CGGACAGATA CAGGTTGAGA
961 CATGTAGAT TTCAATTTAT ATGCTACAGA TATTTGTATA CAGCAGAGC GTGCACAGC
1021 ACAGGTACT TTAAATATG ACAGGCCA AAGGATGCA CAACAGTGCG GTGCCCTAAA
1081 ACGAATGAT ACAGACAGTA TAGAAGCAG CCCTTTAGCA AAGTCCCATC TACACAGCT
1141 ATCAATGAT GTACAGACGA CACAGCAG ACACCGGCG TATACAGTGC CGSACAGCG
1201 CTATGCCAAT ATGCAAGTG AACTAATC GGAAGTAACT GTAGCACTA ATACAAATG
1261 GGGGACGGG GAGATGAG GGAATATG CCACAGATA CGGAGGACT GTAGTAGT
1321 ACAGATGCT ATGATATG AAAACAGGA TCCTAAATCA CTACTACGC AACTAAAAT
1381 ATTTATCAA TGTAAATAA AAAAGCTGC AATGTAAAC GAATTTAAA AGATATAG
1441 ATTTACTCT TATGACTAG TACGTACAT TAAAGTGT ATAGCAATC ATGACGACTG
1501 GTPAGCACA ATATTGGAG TAAATCCAAC CATTCGCCGA GGGTTTAAAA CACTAATTA
1561 ACATATGCA TTATATACC ATATACATG TTATAGTACA AAAACCGAG TATTAATAT
1621 ATGTTTATG AGATACAGA GTGGGAAA TAGAATAA GTAGAAAGC GAATTAATC
1681 ATGTGTGAT GTCCAGACA CGCGTATGCT TTTCGACCA CBAATTTGC GTAGCCCTG
1741 TGCACATG TATTGTRTA GAACAGGAT ATCTAATAT ATGAGAGTGT GTGGACAC
1801 CCAATATGG ATAAAGATG ACAATATG ATACATGAT ATAGATGATA GTGTATTGA
1861 TCATACAG ATGTAGCAT GGGCATTTGA TAATGATTA ACAGATGAA GTGTATTAG
1921 ATTTTCACT GATATGAT CAGATTTGTA TAGTAATGCT GACGCTTTT TAAAGCAAA
1981 CTGTACACA AATATGTA AGATTTGTC AATATGCT AGACATTA AACGGGCAA
2041 AAAACGACA ATGTCAATG CGCAATGAT TAAATTTAGA TGCAGTAAAT GTATACAGG
2101 CGGTATTTG CACCAATTTG TACAGTTTAT CAGCATTTT CAGCATTTAT TTATATACA
2161 TTTATGTCGA TTAAAGAT TTTTAAAGG CACGCCAAA CGTAATTTGA TAGTTATCA
2221 TGGGCGACCA AATACAGCA AGTCATATTT TTGCAATGCT CTATACATC TCTTCGAGG
2281 CACAATAT TTCAATGTA ATTCAGCTAG TCACITTTTG TTAGGCCAC TTGCAATGC
2341 AAAAATAGC ATGTATAGT ACCACAGG TACATGTTG TCAATTTTG ATATTTACAT
2401 GAGAATGCA TTAGATAGA ACCCAATAG TTTAGATAGA AAACACAGC ACCTAATAC
2461 AATAAATGT CACCAATGC TAATAACATC AATATCCAA CCTGTGGAAG ACAATAGGTG
2521 GCCTATTTA CMTAGGAG TAACATGCT TAAATTTCTT AATGCATTC CATTTGACCA
2581 TTAGGTCGAAA TTACAGGAT CACAGTATA TAAAACTG AAATGTTTTT TTGAAAAGC
2641 TTGCGTACAA TTACAGTTC AGCAGACGA GGATGAAGGA GACAAATGAT AAGTGAAT
2701 CCAACGPTT ACAGATTTA CAGGAGAAA TATTAGAAC TAATTAGACA GACATCAAT
2761 GTAFAGAGA CAATTTAC TATGGAATG GTATACGCT GBAATGCA ATATATTATG
2821 CAGCAGAGA ACGTGTGATG CATATATG ACCACAGGT GTAGCTCTT GTAAACATTT
2881 CAAAACATG ACGATATCA GCCATGAAC TCGAGATGG ACTAGAGAGC ATGTCTAAA
2941 CTGCTATAG TCGAGAGAG TGGACATTA GGGACACAAG TAATGAACTA TGGCATACA
3001 AGCCAAGCA ATGTTTAAA AACATAGTG TTACAGTGA AGTATGTTAT CAGGTCAGA
3061 AAGTATAC TCATGATTA ATAGATGAG GTACAAITTA TTTTAAAC AGTACAGACA
3121 CATGTGPTA AACGAAAGG TTTGTGAATG ATTTGGGTGT ATATATATG TATGAAAAC
3181 AAAAAAGCA TTATAGATG ATTTGAATG ATGCACACT ATATGGAAT AGTGAAT
3241 GGGACGTCG TTATAGATG ACATATTC ATTGTCCTGA CTTATNGTC AGTACACTG
3301 AGGACAGCT ACCCATTC GAATCTATG CCGCACTACA GACACACG CGACCCATA
3361 CMACCGANG TGCCCAACG ACCAAAAA CAGCTCGAG GTGTCTTG AAGCACCCA
3421 GAGATAGCG AATCATGAG CCCTCTGAGC CCAAGACGT GTCCGTGAGC TGTGTACCC
3481 TCCACTACT ANTAGAACT GCAGGCCA ACAAAAGG GACGTTGT AGTGTACCA
3541 CTACACCTAT AGTCACTTTA AAAATGGTT TAAATGGTTT AAGTATAGT AGGTATAGT
3601 TGCAAAATA TAAAGCTTTG TATGAAAA TATCATGTAC ATGCACTTG ATACGGGTA

3661 GGGATCAAC CAATACAGGA ATATTGACT TAACATATAG TACTGAAGCA CAACGCCAGA
3721 AATTTTGGGA AACTTTAAA ATACTTCCA GTGTAACTGT TTCACTAGCA TATATGACAT
3781 TATGATGCTT TATTGTACC ACATCTATA CTGTATGATC ATTGGATACA GTACACATA
3841 CATGTACATA TGATTGTAC GTATTTTTT GTGTGTTTGT GTGTATGATC TATATAATG
3901 TGCACATGCC CGCTCTGCA CCAATGCA CTGTGTGCTG ATGTGTGAT ACCTGTGTTT
3961 GTGTTTATAT TAGTACGTAC CACACATG GAGGCTTTG CTGTATATAT ACTTTTITT
4021 TTTACTGCTA TGTGGGTTT ACACGTTT GCTGTTATA GTATGCCTTA AGTTTGTAT
4081 TTGTGATTTG TATTGGTGA TATTTGTATA AATAAATATG GTATCACCG CAGCTGCCAG
4141 GCGAAGCTG GCATCTCAA TGAATTTATA TAAAACATGC AAACAATCAG GCACATGCC
4201 TCTCATGTTT ATAAATAAG TTGAAGGCAC CACATTTGCA GACAAATAT TGCATATGAC
4261 CAGTTTATGTT ATTTTITTTG GTGCTTAGG CATTTGGTAT GAGCTAGGAA CCGGGGTGCG
4321 TACTGGGTAG ATTCTTTTAG TTGTAAACG TAATACTGTT TTAGATGTTT CGCCTGCACG
4381 TCCACCTGTC CTGTGTAACG CTGTGGTCC TACAGAACG TCCATGTGCG AATGTGTGGA
4441 AGATTCCAGT GTATTACAT CTGGCACAC TTTCAGACA TTTCAGGGA CTTCGTGGTT
4501 TGAATATACA TCTTCTCTA CCACATCAC TCTGTGTTA GAATPACCC CTTCGTCTGG
4561 GCTGTGGA CAATGAGTA CTAGTTTATC TAACTCTGCA TTTCAGACC CCAATATAT
4621 AAGAAGTCT CAACAGGTG AAGTCTGTG TAAATGTGTT TTAAGTACCC CCACATGGG
4681 AACATATGGA TATAGAAAA TACTATGCA GTATTTTGA ACACATGCG CAGTACAGA
4741 ACCATATAGT AGTACACTA TACCTGGGT TAGTCGTGTG CGAGGCCAC GTTATATAG
4801 TAGGGACAT CAACAGTTC GTGTAGTAA TTGTGATTTT TTAATCTACC CTTCATCAT
4861 TGAATATCTT GATATCTCTG TGTTCAGCG CTTCGTGACT ACATATCAT ATGAACCTGC
4921 TGACATGCT CTGATCCGG ATTTCTGGA CATTCGTCT TTACATAGC CTGCCTAAC
4981 TTCCGAGA GGCACATG GTTTACAG AGTATGCAA AAGGACATA TGTTTACAG
5041 CCGGGTACA CAATCGGG CACAGTGC CATTTATCAT GATATAGT CTAATGTACC
5101 TCGTCTGAC ATGAACATC AACCTTTGT TGCCTGAGC CAGTCTGACC GTATGTATG
5161 TTTATATGAT ATATACGAC CAGATATGA CAATATACA GTATTGATC TACATATCCA
5221 TAAATGCTACA TTTACTCCC GTTCCCATAT ATCTGTTCTT TCATTAGCT CTACAGATC
5281 TACTACATAT GTACACTA CATTCCTAT TGTACTGCT TGAACACGC CTGTAAATAC
5341 TGTCTCTCAT GTTGTGTAC CAGCAACGTC TCACAGTTC CTTTAAAC CTTCTACAC
5401 ATTTGATACA ACCTATGCCA TACTATATA TGGCAACAAT TATTTATTT TACCATAT
5461 GTTCTTTTAA TTAATAAAC GTAAACGCT TCTTTATTT TTTCGAGATG GCATATGCG
5521 GCTCTAGCA CAACTGGTG TATTGCCCT CCCCCTAGT GCGAGAGTT GATACAGAT
5581 ATGATTAAGT ACACGACT GGCATTAAT ATCTGCTGG TACATCTAGG TTTATACTG
5641 TAGCAATCTC ATATTTTAG GTCCCTATG CTGGGGCGG CAAGCAGAC ATTCTAAGG
5701 TGTCTGATCA TCAATACAG GTGTTTAGA CCGCTTACC TGATCTAAT AATTTATG
5761 TCTCTGATG TACATATAT AACCTGATA CCGACGATT GTTATGGCC TGTGTTGGT
5821 TTGAATATAG TGGGGGAG TCCCTAGGT TTGCGCTTAG TGGGCATCCA TTATATAA
5881 GCTGTAGTAC TACTGGCAAT TCCCTCTTTC AGCTCCAAA AAATCTAAG CACAGTAGG
5941 ACAATTTTTC AGTGACTAT ACAAACGC AACTATGTAT TATAGCTGT GTTCTGTCCA
6001 TTGGGGACA CTGGGCCAA GGAATCTCT ATAGCCTAG CAATGTCAG CCGGGGCT
6061 GTCCACAT TGAATTAGTA BATACACTA TTACAGATG GATATGAT GATPACAGAT
6121 ATGGTCTAT GGCATTAGT ACATACAG AACAAAAAG GAGGTGCTT TTATATAT
6181 GATCTACAGT CTGCAATAT CTGACTAT TCAAAATGC TGCAGATGTA TATGAGACA
6241 GTATGTTCTT TTGTTTACG AGGAACAGT TATTTGCTAG GCAITTTGG AATGAGGG
6301 GCACTGTAGG GGACACTATA CACTATGAT TGTATATAA GGGCACTGAC ATACGTGACA
6361 GTCTGTGAT TTATGTATG CCCCCTCGC CTAGTGGGTG TATGTTATCC TCAGACTCCC
6421 AGTTATTTAA CAAGCCTAT TGCCTGACA AGGCACAGG ACACAAAT GGTATTTGT
6481 GGCATATCA ATTTTCTT ACTGTTGTG ATACCCTCG CAGTACCAAT TTTACTTTGT
6541 CTACTACTAC TGAATCAGT GTACCAATA TTTATGATCC TAATAATTT AAGGAATATA
6601 TTAGCGATG TGAAGATAT GATTTGCAAT TATATTTCA GTTGTGACT ATACATGAT
6661 CCACTGATG ATGTCTCTAT ATACATACTA TGAATCTCG TATTTGGAT TATTTGAAT
6721 TGTGTTTGC CCCCACCA TCTGTAGTC TTGTAGATAG ATACCGTAT CTGCAATCAG
6781 CAGCAATTTAC ATGTCAAAA GAGCCCTCG CACTACTAA AAGGATCCA TATGATGCT
6841 TAAACTTTTG GAATGTAAAT TTAAGAAAA AGTTTATTC TGAACCTGAC CAGTTTCTT
6901 TAGCAGCAGA ATTTCTTTTA CAGCAGCGC TCCGCCAGC ACCCACTATA GCGCCCGTA
6961 AAGCCGCCG CACAGCACT ACTGCACTA CCTCTAAGCA CAACAGTAAA CGTGTGTCAA
7021 AGTATTTGTT GTATGTTTGT TTTGTGATG TATATGTTG TGTGTTGTTA TATGTGTCAT
7081 GTTGTGTTGT CAGTGTGAG TATATGTTG TGTGTATGT TGTGTATGT TGCAGTATG
7141 TTTGTATAT CTGTTTTGT TATATAAT TGTATGTCAG TTTACTTGT GTTGCACCC
7201 TGTGACTAAC ATATGCTCTT TTTTACATA CATAGAGAT GCACATTTT CTATATAT
7261 TGTAGCGCC CCAATAGGTG TGTATAGTA TATATAAT ATATAGTT CTATATATA
7321 CCAAGGCCG ATTTTAAA GCAATTTAG AACATGTTG TTCAAGAAA ACATGTTTCA CTTGGTTTA
7381 TCCTTCTATA CAGTATTTAA AACTATGTT TTAGCAACAAA CATTGTTTCA CTTGGTTTA
7441 CCCACATATG TGGCACCGT AACATATGT ACTAGCGCG CTATCTAGT CATCATCTG
7501 TCCAGGTGCA GTGCAACAT AGTTTGGCAG CCTATATATC TCCACCTTG TAATAAACT
7561 CTTCTTTAGG GTGCACTTAA TACTGTTTGT ACTTGCCTAA TACAGTAGT GGCTGTATA
7621 ACTTACTTGT CATTGTAAA TGTGCTTTGT AGTGTATGT ATACAGTAG TATACAGTA
7681 TCCAATAAT TGTGCATGC AAATAGTTG GCACACATA CCAATCTTT TACTTATAC
7741 ATTTTACAT CATTTTATAG TATAAGGGA AGGTATGCA CAGGTATCA CCGAAAAAG
7801 TGTATATAA GCTGAAACA CAGTTTGTCT ATACCA

HPV 68b-CIN2-DeI

The sequence is described in Results section 2.1.4, including Figures 2.1.5-2.1.6 and Table 2.6.

1 ATGGGCGTAT TTCAACACC TGAGGACGG CCATACAAT TGCAGACCT GTGAGGACA GTGAGGACA
61 TTGGACACCA CACTGCATGA GCTTCAATA GACTGTGCT ATTTCAGAG GGCACTACAA
121 CGGACAGAGG TATATGAAT TTGCTTTGGT GACTTAAATG TAGTATAGG GCAAGGGGTA
181 CATTATAGCT GATGCAATC AATGATTAAT TTTTATCGGA AATATACGGA ACTACGATAT
241 TACTCAGAT CGGTATGTC AACACATTA GAACACATTA CTRATACAA GTTATATGAT
301 TTATCAATPA GTGCGATG TGCGCTGAAA CCATGTAGTC CTGCTGAAA ACTAAGGCAC
361 TTATTTTCAA AACAAGAT TCATARAAT CCGAGAACCT TTACAGACA GTGTCCGAC
421 TGTGTGAGCA GTAAAGAGA GGCACCGAGA CGCAACCGC TTACAGACA GTGTCCGAC
481 AAGATATGAT GGACATGCA CACCGTGC CAAGATTTGT TTAGAGTTAT GTCCATGAA
541 TGAATATGAT CGGTGTGACC TTGTATGTCA CGAGCAATTA GAGATTTAC AGCATGAAT
601 AGATGAAGCC GACCATATC TTATACCA CAACATCAAA CTACTAGCA GAGGGAGCA
661 AATGAAGCTG CACAAATCT AGTATACGT TTGTAACTGT AACAACTAC TGCAACTAGT
721 AGTACAGCG TCAGGAGCA ACTCGGAA CCGTAGCTG CTGTATTAG ACTCACTAA
781 TTCTGTGCT CGGTGTGTC AAGGGAAC CCAGTATCT GCATGCGCA ATTTGAGAG
841 TACAGATGG GACGGAGCG GTGTAAACG ATGTTTTT GTACAAGCA TAGTAGATA
901 ACAAACAGGT GACATGCT CAGAGATGA GATGAAMC CGGACAGCA CAGTTTCA
961 CATGTAGAT TTCAATGCT ATGCTACAGA TATTGTATA CAGGACAGC GTGACAGCG
1021 ACAGGTAGT TTAAATATG AACAGCCCA AAGGATGCA CAACAGTGC GTGCCCTAAA
1081 ACAGATGAT ACAGACATGA TAGAAGCAG CCCTTTAGCA AAGTCGCCAT TACAGAACT
1141 ATCAATAAAT GTAGACATGA CACAGCAG ACACCGGCG TATACAGTGC CGGACAGCG
1201 CTATGCAAT ATGACATGG AAACTAACCT GGAGTAACT GTAGCAACTA ATACAAGCG
1261 GCGCAGCGG GAGATGAAG GGAATAATGG CGACAGCAT CCGGAGACT GTAGTAGTG
1321 AGACAGTCT APACATGAT GAACACAGA TCCTAATCA CCTACTACG AACTAAATG
1381 ATAAAACCTG GAAATGTTT TTTGAAAAG CTTGTGCAA ATTAGACTTG CAGCAGAGC
1441 AGATGAGG AGACAATGAT GAACACTGT TCCACAGCT TAAATGTGT ACAGAGAA
1501 ATATAGAAC ATTTAGACA GGACGTAAA TGTATAAGG ACCATATAA CTATTGAAAC
1561 TGTATACAG TGAATAATG ATATATAT TCAAAATCT AAGCATATCA AGCATATTA
1621 GACACACAG TGGTGCTCC TGTATACAT TTCAAAACTA AAGCATATCA AGCATATTA
1681 CTGACAGAT CACTGATGCA CATGTCTCAA ACTGCATATA GTGCAGAGA GTGGACATTA
1741 AGGACACAA GTPATGACT ATGGCATCA AAGCAAAGC AATGTTTTAA AAACATGTT
1801 GTACAGTGG AAGTATGAT TGAGCTGAC AAAGCAACT CAATGCATCA TGTAGTAGG
1861 GTTACATTT ATTTTAAAA CAGTACAG ACATGCTGA AAGCGAGG GTGCTGGAT
1921 TATTGGGTG TATATATAT GTATGAAA GTATGAAA ATTAAGAG GTTTATGAT
1981 GATGCACAC ACTATATG TGTGAAAA TGGGACGTG ATATATAGG CAACTAAT
2041 CATTGTCTG ACTATATG ATGATATG CAGCAACAG TATCCACTC TGAATCTATT
2101 GCGACCTCA GATGCTCTG ACCACCCAG GTGCCCATC GTGCCCATC CACCAAAAA
2161 ACACGCTCA CCGCTCTTG CAGACCCAG AGACAGTAG GAATCACTGA GCCCTCTAG
2221 CCCACGACG TGTCCGTGA CGGTGTCAAC CTCCCACTCC TCAGTAGAG TGCAAGGCC
2281 ACAAAGAA GGAAGTGG TTGTTGTGAC ACTCACCTA TAGTGCATT AAAGGTGAC
2341 AAAATGCTT TAAATGCTT TAGTATAG TTGCATAAG ATAACGCTTT GTATGAAT
2401 ATATCATGTA CATTGCAATG GATAGGGGT AGGGATCAA CCAATACAG AATATGACT
2461 GTACATATA GTACTGAG ACACGCCAG AAATTTTGG TGTATGTAC CACACTGTAT
2521 AGTGTACTG TTTCAATAG ATATATGCA TTTATAGTGT TGTATGTAC CACACTGTAT
2581 ACTGTATGA TATTGGATG ATACACCAT ACATGTATC ATGATGTAC TGTATTTTT
2641 GGTGTGTTT TGTGTGGA TGTATATG TGTGACTGC CCGCTCTGC AGTCATGTA
2701 TGTGTGTTG TGTGTGGA TACTGTGTT TGTGTTTATA TACTGTGTA CCACACCAT
2761 GGAGTCTTT TGTGTATATA TACTTTTTT TTTTACTGCT ATGTGGGTAT TACACATTT
2821 TGCTCGTTT AGTATGCTT AAGTTTGTTA TTTTGCATTT GTATTGGTAT ATATTGTAT
2881 AATAAATAT GTATACAC CGTGTGCA GGCACAGCG TGCATGCA ACTGAATAT
2941 ATAAACATG CAACATCA GGCATATC CTCTGATGT TATAAATAG GTTGAAGCA
3001 CACACTTGC AGACAATA TTGCATGGA CAGTTTAGT TATTTTTG GTGGCTTAG
3061 GCAATGTPAC TGGGTAGA AGCGGGGTC GTACTGGTA CATTCCTTAA GTGTGAAC
3121 CTAATACAT TGTATAGT TCGCTGCA GTCCACCTGT GGTATTGA CCGTGGGTC
3181 CTACAGAAC CTCATTTGT CAATGTGTG AAGATTCAG TGTATTACA TCTGGCAC
3241 CGGTACAAAC ATTTACAGC ACTCTGGGT TTGAATATAC ATCTCTCT ACCATCAC
3301 CTGCTGTTT AGACATPAC CCGTCTGTC GGTCTGTGCA AGTAAGCAGT ACTAGTTTA
3361 CTACCTTGC ATTTGACAG CCGCTATTA TAGAAGTGC TCAACAGGT GAAGTCTCT
3421 GTATGTGTT TGTAAATG TCCACATCG GAACATGAG ATATGAAGA ATACATGC
3481 AGTATTTGC ACACATGCG ACTGTACAG AACCTATTAG TAGTACCT ATACCTGGG
3541 TTATGCTGTG GCGAGGCA CTTTATATA TTGTATCAT TGTAACTCT GGTGTAGA
3601 ATTTGTAFTT TGTAACTAC CTTTATAT TGTAACTAT TGTAACTCT GGTGTAGC
3661 CTGTTGTAFT TACACTAC TGTAACTG CTGACATAG TGTAACTCT GGTGTAGC
3721 ACATTTGTC TTATACATG CTTCCGAAG AGGCACAT CTTTGTAGCA
3781 GAGTAGCAA AAAAGCAAC TATGTTTACA CGCCGGGTA CACAAATCG GGCACAGTG

A3. Partial URR sequences of eleven HPV68 positive samples

The nucleotide sequences of eleven clinical samples (Results, Table 2.9) covering partial HPV68 URR region, corresponding to pos. 7279-7769 of the HPV68b-CIN2, are shown as an alignment. The sequences are described in Results section 2.1.5, including Figure 2.20 and Table 2.10.

Reims-06	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	300
Reims-09	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
Reims-12	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
Reims-14	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
Reims-15	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
Reims-16	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
Reims-18	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
Reims-28	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
Reims-31	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
Reims-33	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
HPV68b-CIN2	ATAGTTTGGCAGCCTATATATCTCTCCACCCTTGTAAATAAACTGCTCTTTTAGGCATAGTTTT	
Reims-06	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	360
Reims-09	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
Reims-12	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
Reims-14	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
Reims-15	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
Reims-16	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
Reims-18	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
Reims-28	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
Reims-31	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
Reims-33	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
HPV68b-CIN2	TTAAGTCTGTTTTTACTCTGCTTAATAGCATAGTTGGCCGTGTATATTAACACTCTTTTGCATTCGAG	
Reims-06	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	420
Reims-09	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
Reims-12	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
Reims-14	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
Reims-15	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
Reims-16	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
Reims-18	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
Reims-28	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
Reims-31	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
Reims-33	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
HPV68b-CIN2	AATCTGTCTAGTAGTGTAAAGTTTATACAGTGAAGTAAATACCAATCCATAAAATTTGTGCAAC	
Reims-06	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	480
Reims-09	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
Reims-12	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
Reims-14	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
Reims-15	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
Reims-16	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
Reims-18	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
Reims-28	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
Reims-31	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
Reims-33	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
HPV68b-CIN2	CGAAATAGGTTGGGCACACATACCAATACCTTTTACTTTTAAACATTTTAAACATCAATTTTAT	
Reims-06	ATAATAAAGGG	491
Reims-09	ATAATAAAGGG	
Reims-12	ATAATAAAGGG	
Reims-14	AGTATAAAGGG	
Reims-15	AGTATAAAGGG	
Reims-16	AGTATAAAGGG	
Reims-18	ATTAAAAAGGG	
Reims-28	ATTAAAAAGGG	
Reims-31	ATTAAAAAGGG	
Reims-33	AGTATAAAGGG	
HPV68b-CIN2	AGTATAAAGGG	

A4. Nucleotide sequences of identified HPV16 integration junctions in ASP16 experiments

Sample HSIL-61979

See Results section 2.2.2.6, pages 78-80. The nucleotide sequence shown here is composed of the junction-specific PCR product sequence and the ASP16 E26 sequence reads (Results, Figure 2.44). Red color indicates HPV16 DNA between pos. 2403-2516, and black color the flanking cellular DNA.

```
1  AGCAGATGCC AAAAAAGGTA TGTGGATGA TGTACAGTG CCTGTGTGA ACTATATAGA
61  TGACAAATTA AGAAATGCAT TGGATGGAAA TTATGTTTCT ATGATGTGTA AGCAGATGTG
121 GGCCAGCAT GCTTACAGCG GGGACACCAC TTCCCAACAG CCTGAAGAAC TCCCATTTCA
181 CCTTAAAAAC
```

Sample HSIL-75857

See Results section 2.2.2.6, pages 81-82. The nucleotide sequence shown here is assembled from the sequences of two junction-specific PCR products (Results, Figure 2.46). Red color indicates HPV16 DNA between pos. 1064-1149, and black color the flanking cellular DNA.

```
1  CACATGCGTT GTTTACTGCA CAGGAGCAA AACACATAG AGATGCAGTA CAGGTTCTAA
61  AACGAAGTA TTTCGGTAGT CCACTTGTGT TTGTAAGTTC TCATCTTAAT TTGCAGTAT
121 CCAGTTTGG AATGCCACGA AAAATAAATT ATGTTTATTA AAATTATTTT TTAAAAAAC
181 CTTTCCAAAT AATTGAAAAA CACTATTCTA CTCAAATCTC AAAAATATCC AAATATACGT
241 TTTAAACAGA AAAATTAAAA AAATATTTTC TCTGGATTCT CTATTTCCTG AACCTTTAAA
301 AAATGTTTTT CAGTGAATAA TATGTTCCGA GGTGATTCCG ACTACAAAT
```

Sample CIN2/3-1801

See Results section 2.2.2.6, pages 82-83. The nucleotide sequence shown here belongs to the junction-specific PCR product of primer pair 1801-2 (Results, Figure 2.47). Red color indicates HPV16 DNA between pos. 1656-1913, green color the 4 bp of unidentified origin, and black color the flanking cellular DNA.

```
1  TCAAAAGTTTA GCATGTTTCAT GGGGAATGGT TGTGTTACTA TTAGTAAGAT ATAAATGTGG
61  AAAAAATAGA GAAACAATTG AAAAATTTGCT GTCTAAACTA TTATGTGTGT CTCCAATGTG
121 TATGATGATA GAGCTCCAA AATTGGGTAG TACAGCAGCA GCATTATATT GGTATAAAAC
181 AGGTATATCA AATATTAGTG AAGTGATGAG AGACACGCCA GAATGGATAC AAAGACAAAC
241 AGTATTACAA CATAGTTTGG CATGCCACT CCATGCCATC AGCCCCAAG TCATGACGAC
301 TACAGGTATC TCCAGACATC ACCAGTGTG CCTGGGGCG GGTGGGGGG CAGGCTTACC
361 CCTTTGAGA ACCTCCTGTC TTGAGGCAG CAAAC
```


1289
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
MRI-H186
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
MRI-H186
SiHa
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
Caski
CIN2/3-0001
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-0002
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-0005
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-1503
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-1511
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-1801
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-2219
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-2227
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-2229
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-2237
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-3009
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-3035
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-4242
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CA-07C381
HSIL-75857
CIN2/3-0004
AAAAAGGAGATTATTTGAAAGCGAAGACAGCGGGTATGCCAATCTGAAGTGGAAACTCA
CIN2/3-4238a
CA-07C368
HSIL-66019
HSIL-61979
LSIL-75022

1349
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
HPV16R
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
MRI-H186
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
MRI-H196
SiHa
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
Caski
CIN2/3-0001
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-0002
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-0005
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-1503
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-1511
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-1801
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-2219
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-2227
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-2229
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-2237
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-3009
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-3035
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-4242
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-4242
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CA-07C381
HSIL-75857
CIN2/3-0004
GCAGATGTTACAGGTAGAAGGGGCCCATGAGACTGAAACACCATGTAGTCAGTATAGTGG
CIN2/3-4238a
CA-07C368
HSIL-66019
HSIL-61979
LSIL-75022

1469
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
HPV16R
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
MRI-H186
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
MRI-H196
SiHa
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
Caski
CIN2/3-0001
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-0002
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-0005
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-1503
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-1511
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-1801
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-2219
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-2227
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-2229
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-2237
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-3009
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-3035
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-4242
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-4242
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CA-07C381
HSIL-75857
CIN2/3-0004
AAGACACACTATATGCCCAACACACACTTCAAAATATTTTAAATGTACTTAAACCTAGTAA
CIN2/3-4238a
CA-07C368
HSIL-66019
HSIL-61979
LSIL-75022

[illegible][illegible]

HPV16R
 MRL1-H186
 MRL1-H196
 SiHa
 Caski
 CIN2/3-0001
 CIN2/3-0002
 CIN2/3-0005
 CIN2/3-0002
 CIN2/3-1503
 CIN2/3-1511
 CIN2/3-1801
 CIN2/3-2219
 CIN2/3-2221
 CIN2/3-2229
 CIN2/3-2237
 CIN2/3-2237
 CIN2/3-3009
 CIN2/3-3035
 CIN2/3-4242
 CA-07C381
 CA-07C381
 HSIL-715857
 CIN2/3-0004
 CIN2/3-4238a
 CA-07C368
 HSIL-66019
 HSIL-619179
 HSIL-715020
 LSiL-6922

[illegible]

HPV16R	CaskI
MRI-H186	TIN2/3-0001
MRI-H196	TIN2/3-0002
SiHa	TIN2/3-0005
	TIN2/3-1503
	TIN2/3-1511
	TIN2/3-1801
	TIN2/3-2219
	TIN2/3-2227
	TIN2/3-2229
	TIN2/3-2237
	TIN2/3-3009
	TIN2/3-3035
	TIN2/3-4242
	CA-07C381
	HSII-75857
	TIN2/3-0004
	TIN2/3-4238A
	CA-07C368
	HSII-66019
	HSII-61979
	L5II-75022

[illegible]

HPV16	HPV18	HPV31	HPV33	HPV35	HPV39	HPV45	HPV52	HPV58	HPV59	HPV68	HPV70	HPV74	HPV82	HPV84	HPV89	HPV91	HPV92	HPV93	HPV94	HPV95	HPV97	HPV98	HPV99	HPV100	HPV101	HPV102	HPV103	HPV104	HPV105	HPV106	HPV107	HPV108	HPV109	HPV110	HPV111	HPV112	HPV113	HPV114	HPV115	HPV116	HPV117	HPV118	HPV119	HPV120	HPV121	HPV122	HPV123	HPV124	HPV125	HPV126	HPV127	HPV128	HPV129	HPV130	HPV131	HPV132	HPV133	HPV134	HPV135	HPV136	HPV137	HPV138	HPV139	HPV140	HPV141	HPV142	HPV143	HPV144	HPV145	HPV146	HPV147	HPV148	HPV149	HPV150	HPV151	HPV152	HPV153	HPV154	HPV155	HPV156	HPV157	HPV158	HPV159	HPV160	HPV161	HPV162	HPV163	HPV164	HPV165	HPV166	HPV167	HPV168	HPV169	HPV170	HPV171	HPV172	HPV173	HPV174	HPV175	HPV176	HPV177	HPV178	HPV179	HPV180	HPV181	HPV182	HPV183	HPV184	HPV185	HPV186	HPV187	HPV188	HPV189	HPV190	HPV191	HPV192	HPV193	HPV194	HPV195	HPV196	HPV197	HPV198	HPV199	HPV200	HPV201	HPV202	HPV203	HPV204	HPV205	HPV206	HPV207	HPV208	HPV209	HPV210	HPV211	HPV212	HPV213	HPV214	HPV215	HPV216	HPV217	HPV218	HPV219	HPV220	HPV221	HPV222	HPV223	HPV224	HPV225	HPV226	HPV227	HPV228	HPV229	HPV230	HPV231	HPV232	HPV233	HPV234	HPV235	HPV236	HPV237	HPV238	HPV239	HPV240	HPV241	HPV242	HPV243	HPV244	HPV245	HPV246	HPV247	HPV248	HPV249	HPV250	HPV251	HPV252	HPV253	HPV254	HPV255	HPV256	HPV257	HPV258	HPV259	HPV260	HPV261	HPV262	HPV263	HPV264	HPV265	HPV266	HPV267	HPV268	HPV269	HPV270	HPV271	HPV272	HPV273	HPV274	HPV275	HPV276	HPV277	HPV278	HPV279	HPV280	HPV281	HPV282	HPV283	HPV284	HPV285	HPV286	HPV287	HPV288	HPV289	HPV290	HPV291	HPV292	HPV293	HPV294	HPV295	HPV296	HPV297	HPV298	HPV299	HPV300	HPV301	HPV302	HPV303	HPV304	HPV305	HPV306	HPV307	HPV308	HPV309	HPV310	HPV311	HPV312	HPV313	HPV314	HPV315	HPV316	HPV317	HPV318	HPV319	HPV320	HPV321	HPV322	HPV323	HPV324	HPV325	HPV326	HPV327	HPV328	HPV329	HPV330	HPV331	HPV332	HPV333	HPV334	HPV335	HPV336	HPV337	HPV338	HPV339	HPV340	HPV341	HPV342	HPV343	HPV344	HPV345	HPV346	HPV347	HPV348	HPV349	HPV350	HPV351	HPV352	HPV353	HPV354	HPV355	HPV356	HPV357	HPV358	HPV359	HPV360	HPV361	HPV362	HPV363	HPV364	HPV365	HPV366	HPV367	HPV368	HPV369	HPV370	HPV371	HPV372	HPV373	HPV374	HPV375	HPV376	HPV377	HPV378	HPV379	HPV380	HPV381	HPV382	HPV383	HPV384	HPV385	HPV386	HPV387	HPV388	HPV389	HPV390	HPV391	HPV392	HPV393	HPV394	HPV395	HPV396	HPV397	HPV398	HPV399	HPV400	HPV401	HPV402	HPV403	HPV404	HPV405	HPV406	HPV407	HPV408	HPV409	HPV410	HPV411	HPV412	HPV413	HPV414	HPV415	HPV416	HPV417	HPV418	HPV419	HPV420	HPV421	HPV422	HPV423	HPV424	HPV425	HPV426	HPV427	HPV428	HPV429	HPV430	HPV431	HPV432	HPV433	HPV434	HPV435	HPV436	HPV437	HPV438	HPV439	HPV440	HPV441	HPV442	HPV443	HPV444	HPV445	HPV446	HPV447	HPV448	HPV449</
-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	----------

[illegible]

HPV16R		Cask1	CIN2/3-0001		CSL1-75857	CIN2/3-0004
MRI-H186			CIN2/3-0002		CA-07C384	CA-07C384
MRI-H196		SiHa	CIN2/3-0005		HSII-66019	HSII-66019
			CIN2/3-1503		CSL1-75022	CSL1-75022
			CIN2/3-1511		CA-07C381	CA-07C381
			CIN2/3-1801		HSII-75857	HSII-75857
			CIN2/3-2219			
			CIN2/3-2227			
			CIN2/3-2229			
			CIN2/3-2237			
			CIN2/3-3009			
			CIN2/3-3035			
			CIN2/3-4242			
			CSA-07C381			

[illegible]

HPV16	CIN2/3-0001	Caski
MRI-H186	CIN2/3-0005	SiHa
MRI-H196	CIN2/3-1503	SiHa
	CIN2/3-1511	Caski
	CIN2/3-1801	
	CIN2/3-2219	
	CIN2/3-2227	
	CIN2/3-2229	
	CIN2/3-2237	
	CIN2/3-3009	
	CIN2/3-3035	
	CIN2/3-3242	
	CIN2/3-4242	
	CA-07C381	
	HSII-73857	
	HSII-73904	
	IN2/3-4238a	
	CA-07C368	
	HSII-66019	
	HSII-61979	
	LSII-75022	

[illegible]

HPV1 6R
 MRI-1H186
 MRI-1H196
 SiHa
 CasKi
 C1N2/3-0001
 C1N2/3-0002
 C1N2/3-0005
 C1N2/3-1503
 C1N2/3-1511
 C1N2/3-1801
 C1N2/3-2219
 C1N2/3-2227
 C1N2/3-2229
 C1N2/3-2237
 C1N2/3-3009
 C1N2/3-3035
 C1N2/3-4242
 CA-07C381
 HPV1 75857
 C1N2/3-0004
 C1N2/3-4238a
 CA-07C368
 HSII1-66019
 HSII1-61979
 HSII1-75022

2009

HPV16R	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
MR1-H186	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
MR1-H196	ATGGGCGCTACGA-----TAGACGATAGTGAATTCGCATATAAATATGACAAATTT
SiHa	ATGGGCGCTACGAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CaSk1	ATGGGCGC-----TAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-0001	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-0002	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-0005	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-1503	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-1511	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-1801	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-2217	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-2219	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-2229	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-2237	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-3009	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-3035	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-4242	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CA-07C381	ATGGGCGC-----TAGACGATAGTGAATTCGCATATAAATATGACAAATTT
HSIL1-75837	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CIN2/3-0004	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CA-4238a	ATGGGCGC-----TAGACGATAGTGAATTCGCATATAAATATGACAAATTT
CA-07C368	ATGGGCGCTACGATAATGACATAGTAGACGATAGTGAATTCGCATATAAATATGACAAATTT
HSIL1-66019	ATGGGCGC-----TAGACGATAGTGAATTCGCATATAAATATGACAAATTT
HSIL1-61979	ATGGGCGC-----TAGACGATAGTGAATTCGCATATAAATATGACAAATTT
HSIL1-75022	ATGGGCGC-----TAGACGATAGTGAATTCGCATATAAATATGACAAATTT

206

HPV18R	GGCAGACCTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
MR1-H186	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
MR1-H196	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
SilHa	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CaSk1	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-0001	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-0002	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-0005	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-1503	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-1511	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-1801	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-2219	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-2227	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-2239	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-2237	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-3003	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-3035	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-4242	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
Ca-07C381	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
HS11-75857	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-0004	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
CIN2/3-4238A	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
Ca-07C368	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
HS11-66019	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
HS11-61979	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT
HS11-75022	GGCGACACTAATAGTAATGCAAGTGCCTTCTTAAAAAGTAATTCACAGCGCAAAAATTTGT

[illegible][illegible]

[illegible]

HPV1 6R
 MRI-H186
 MRI-H196
 SilHa
 CasKi
 C1N2/3-0001
 C1N2/3-0002
 C1N2/3-0005
 C1N2/3-1503
 C1N2/3-1511
 C1N2/3-1801
 C1N2/3-2219
 C1N2/3-2227
 C1N2/3-2229
 C1N2/3-2237
 C1N2/3-3009
 C1N2/3-3035
 C1N2/3-4242
 C4-07C381
 HSIL-75857
 C1N2/3-0004
 C1N2/3-4238A
 C4-07C368
 HSIL-66019
 HSIL-61979
 HSIL-75022

[illegible]

254

HPV18	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
MR1-H186	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
MR1-H196	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
SiHa	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CaSk1	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-0001	-----TTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-0002	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-0005	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-1503	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-1511	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-1801	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-2219	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-2227	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-2229	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-2237	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-3009	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-3035	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-4242	AAATTTAGTTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
Ca-07C381	AAA-----TTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
HSIL-75837	-----TTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-0004	-----TTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
CIN2/3-4238A	AAA-----TTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
Ca-07C368	-----TTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
HSIL-66019	-----TTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
HSIL-61979	-----TTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT
HSIL-75022	-----TTTCTATGATGTTAAAGCATAGACCATGGTGTACAACTAAATAGCCCTCCATT

[illegible]

[illegible][illegible]

[illegible][illegible]

A6. Nucleotide sequences of HPV16 ORF E6 of 25 DNA samples in ASP16 experiments

See Results section 2.2.2.9. The complete HPV16 ORF E6 nucleotide sequences of the 25 DNA samples (see Results, Table 2.13 and Table 2.18), covering pos. 83-559 are shown as an alignment.

HPV16R	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	262
CIN2/3-0004	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CA-07C368	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
MRI-H186	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-2229	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-4238a	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
HSIL-66019	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
HSIL-61979	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-0001	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-0002	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-1503	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-1801	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-2227	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
HSIL-75857	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
MRI-H196	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
SiHa	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
Caski	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CA-07C381	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
LSIL-75022	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-3009	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-3035	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-4242	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-0005	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-1511	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-2219	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
CIN2/3-2237	ATGCACCAAAAGAGAACTGCAATGTTTCAGGACCACACAGGAGCCAGCCAGAAAGTTACCA	
HPV16R	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	202
CIN2/3-0004	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CA-07C368	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
MRI-H186	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-2229	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-4238a	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
HSIL-66019	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
HSIL-61979	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-0001	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-0002	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-1503	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-1801	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-2227	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
HSIL-75857	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
MRI-H196	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
SiHa	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
Caski	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CA-07C381	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
LSIL-75022	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-3009	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-3035	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-4242	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-0005	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-1511	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-2219	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	
CIN2/3-2237	CAGTTATGCAACAGAGCTGCAACAACTATACATGATATATATTTAGAAATGTTGTACTGC	

HPV16R	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	262
CIN2/3-0004	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CA-07C368	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
MRI-H186	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-2229	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-4238a	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
HSIL-66019	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
HSIL-61979	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-0001	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-0002	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-1503	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-1801	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-2227	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
HSIL-75857	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
MRI-H196	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
SiHa	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
Caski	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CA-07C381	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
LSIL-75022	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-3009	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-3035	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-4242	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-0005	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-1511	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-2219	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
CIN2/3-2237	AAGCACACAGTTACTGCGACGTCGAGGTATATGACTTTTCTTCGGGATTTATGCATAGTA	
HPV16R	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	322
CIN2/3-0004	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CA-07C368	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
MRI-H186	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-2229	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-4238a	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
HSIL-66019	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
HSIL-61979	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-0001	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-0002	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-1503	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-1801	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-2227	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
HSIL-75857	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
MRI-H196	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
SiHa	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
Caski	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CA-07C381	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
LSIL-75022	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-3009	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-3035	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-4242	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-0005	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-1511	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-2219	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	
CIN2/3-2237	TATAGAGATGGGAATCCATATCGTGTATGTGATATAAATGTTTAAAGTTTTATTCCTAAAAAT	

HPV16R		GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATATAGGGGTGGTGGACC	559
CIN2/3-0004	CA-07C31	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CA-07C316	LS11-75022	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
MRI-H186	CIN2/3-3009	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-2238a	CIN2/3-3035	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-4238a	CIN2/3-4242	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
HS11-66019	CIN2/3-0001	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
HS11-61979	CIN2/3-0002	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-0001	CIN2/3-1503	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-1503	CIN2/3-1801	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-2227	CIN2/3-2227	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
HS11-75857	HS11-75857	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
MRI-H196	MRI-H196	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
SiHa	SiHa	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
Caski	Caski	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CA-07C381	CA-07C381	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
LS11-75022	LS11-75022	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-3009	CIN2/3-3009	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-3035	CIN2/3-3035	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-4242	CIN2/3-4242	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-0005	CIN2/3-0005	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-1511	CIN2/3-1511	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-2219	CIN2/3-2219	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
CIN2/3-2237	CIN2/3-2237	GAAGAAGACATCTGGCAAAAACGAAAGATTCATATATAGGGGTGGTGGACC	
HPV16R		GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	559
CIN2/3-0004	CA-07C31	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CA-07C316	LS11-75022	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
MRI-H186	CIN2/3-3009	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-2229	CIN2/3-2229	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-4238a	CIN2/3-4238a	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
HS11-66019	HS11-66019	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
HS11-61979	HS11-61979	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-0001	CIN2/3-0001	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-0002	CIN2/3-1503	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-1503	CIN2/3-1801	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-1801	CIN2/3-2227	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-2227	CIN2/3-2227	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
HS11-75857	HS11-75857	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
MRI-H196	MRI-H196	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
SiHa	SiHa	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
Caski	Caski	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CA-07C381	CA-07C381	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
LS11-75022	LS11-75022	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-3009	CIN2/3-3009	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-3035	CIN2/3-3035	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-4242	CIN2/3-4242	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-0005	CIN2/3-0005	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-1511	CIN2/3-1511	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-2219	CIN2/3-2219	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	
CIN2/3-2237	CIN2/3-2237	GTCGATGATGTCCTTTGTCAGATCATCAAGACACGTTAGAGAAACCCAGCTGTAA	

[illegible]

A7. Contents of the file required for ASP16 data analysis

See Results section 2.2.1.3. The contents of a text file, named “hpv16R_EEprimer_151nt”, are shown below.

[illegible]

A8. Source codes of ASP16 data analysis computer programs

See Results section 2.2.1.2 and 2.2.1.4. The Perl source codes of the computer programs for the ASP16 data analysis are shown. Four source codes for the four main program sets are described first, followed by the source codes of thirteen sub-programs for program sets 2-4.

Program set 1:

File name: **PROGRAM_SET_1_mac.pl**

Source code:

```
#!/usr/bin/perl

## About the script #####
# Created by Sasithorn Choteutmontri, Jan 2009.
#####

##### MAIN PROGRAM BODY #####
#####

use strict;
use warnings;

#####
##### PLEASE CHANGE THE PATHS FOR PROGRAMS LOCATIONS AND RESULTS LOCATION #####
#####

##### 1 ##### Path where all Perl (scripts) programs are located on your computer:

my $Pathprog = "/Volumes/Ma-ch02/DKF-Z/PERL/";

#####
my $resultPath = "/Volumes/Ma-ch02/DKF-Z/";

#####
##### USERS area ENDS #####

my $PREFIX_path = $resultPath;
my @data1 = @ARGV;
my @STACK = ();
my @data1_name = ();
my @data1_seq = ();

my $pyroNum = int($InputFile[1]);
my $path = $PREFIX_path."Pyro".$pyroNum."_result/";

system ("mkdir", $path);

my $ba01 = 'TGAC'; # for sample 1
my $ba02 = 'AGAC'; # for sample 2
my $ba03 = 'TCAC'; # for sample 3
my $ba04 = 'ACAC'; # for sample 4
my $ba05 = 'TGTC'; # for sample 5
my $ba06 = 'AGTC'; # for sample 6
my $ba07 = 'TCTC'; # for sample 7
my $ba08 = 'ACTC'; # for sample 8
my $ba09 = 'CTGA'; # for sample 9
my $ba10 = 'CAGA'; # for sample 10
my $ba11 = 'CTCA'; # for sample 11
my $ba12 = 'CACA'; # for sample 12
my $ba13 = 'ATGC'; # for sample 13
my $ba14 = 'ATGC'; # for sample 14
```

```
my $ba15 = 'TACA'; # for sample 15
my $ba16 = 'ATCA'; # for sample 16
my $ba17 = 'TGCA'; # for sample 17
my $ba18 = 'TCGA'; # for sample 18
my $ba19 = 'AGCA'; # for sample 19
my $ba20 = 'ACGA'; # for sample 20
my $ba21 = 'CATC'; # for sample 21
my $ba22 = 'CTAC'; # for sample 22
my $ba23 = 'CTGC'; # for sample 23
my $ba24 = 'CAGC'; # for sample 24

my @barcode = ( $ba01, $ba02, $ba03, $ba04, $ba05, $ba06, $ba07, $ba08,
$ba09, $ba10, $ba11, $ba12, $ba13, $ba14, $ba15, $ba16,
$ba17, $ba18, $ba19, $ba20, $ba21, $ba22, $ba23, $ba24 );

my @smp_le01_name = (); my @smp_le01_seq = ();
my @smp_le02_name = (); my @smp_le02_seq = ();
my @smp_le03_name = (); my @smp_le03_seq = ();
my @smp_le04_name = (); my @smp_le04_seq = ();
my @smp_le05_name = (); my @smp_le05_seq = ();
my @smp_le06_name = (); my @smp_le06_seq = ();
my @smp_le07_name = (); my @smp_le07_seq = ();
my @smp_le08_name = (); my @smp_le08_seq = ();
my @smp_le09_name = (); my @smp_le09_seq = ();
my @smp_le10_name = (); my @smp_le10_seq = ();
my @smp_le11_name = (); my @smp_le11_seq = ();
my @smp_le12_name = (); my @smp_le12_seq = ();
my @smp_le13_name = (); my @smp_le13_seq = ();
my @smp_le14_name = (); my @smp_le14_seq = ();
my @smp_le15_name = (); my @smp_le15_seq = ();
my @smp_le16_name = (); my @smp_le16_seq = ();
my @smp_le17_name = (); my @smp_le17_seq = ();
my @smp_le18_name = (); my @smp_le18_seq = ();
my @smp_le19_name = (); my @smp_le19_seq = ();
my @smp_le20_name = (); my @smp_le20_seq = ();
my @smp_le21_name = (); my @smp_le21_seq = ();
my @smp_le22_name = (); my @smp_le22_seq = ();
my @smp_le23_name = (); my @smp_le23_seq = ();
my @smp_le24_name = (); my @smp_le24_seq = ();

@data1 = getFileName ($InputFile[0]);

my $barcode_sample = "barcode-to-sample-arrangement.txt";

print "\n\n=====
print "\n\nPROGRAM 1 (SORTING) STARTS.....\n\n";

open (BARSAM, ">$path$barcode_sample");
print BARSAM "\n\n";
print BARSAM "\nDefault values for barcode-sample pair are used.\n\n";

print BARSAM $ba01." for sample 1\n"
$ba02." for sample 2\n"
$ba03." for sample 3\n"
$ba04." for sample 4\n"
$ba05." for sample 5\n"
$ba06." for sample 6\n"
$ba07." for sample 7\n"
$ba08." for sample 8\n"
$ba09." for sample 9\n"
$ba10." for sample 10\n"
$ba11." for sample 11\n"
$ba12." for sample 12\n"
$ba13." for sample 13\n"
$ba14." for sample 14\n"
$ba15." for sample 15\n"
$ba16." for sample 16\n"
$ba17." for sample 17\n"
$ba18." for sample 18\n"
$ba19." for sample 19\n"
$ba20." for sample 20\n"
$ba21." for sample 21\n"
$ba22." for sample 22\n"
$ba23." for sample 23\n"
$ba24." for sample 24\n\n";
```

```

sub printData {
    my (@nameqq) = @_;
    use strict;
    use warnings;

    print @nameqq;
    print "\n\n";
}

sub getFileData {
    my ($filename) = @_;
    use strict;
    use warnings;

    my @filedata = ();
    unless ( open(FILE_DATA, $filename) ) {
        print STDERR "\n\nThe program can't open the files \"$filename\"\n\n";
        print "Please re-check input file name and its location\n";
        print "The correct command should be : \n\n";
        print "\t\tperl PROGRAM(location-name) INPUTFastAFILE(location+name)\n\n";
        print "The command should be given under directory";
        print " c:/Perl/bin> in case of Dos Terminal\n\n";
        exit;
    }
    @filedata = <FILE_DATA>;
    close FILE_DATA;
    return @filedata;
}

## getting a STACK array , membership# = line# of original file data
# 0 = name
# 1 = sequence
# 2 = others
sub extractFastStack {
    my (@fastafileData) = @_;
    use strict;
    use warnings;

    my $stackcount = scalar @fastafileData;
    my @stacklist = ();
    for (my $m = 0; $m < $stackcount ; $m++) {
        if ($fastafileData[$m] =~ /\s*$/) {
            push @stacklist, 2;
        } elsif ($fastafileData[$m] =~ /\s*$/) {
            push @stacklist, 2;
        } elsif ($fastafileData[$m] =~ /\s*$/) {
            push @stacklist, 0;
        } else {
            push @stacklist, 1;
        }
    }
    return @stacklist;
}

## get the sequences into an array
sub extractFastSeq {
    my ($data1, $STACK, $data1_seq) = @_;
    use strict;
    use warnings;

    my $stackcount = scalar @data1;
    for (my $p = 0; $p < $stackcount ; $p++) {
        if ($$STACK[$p] =~ "0") {
            next;
        } elsif ( ($$STACK[$p] =~ "1") && ($$STACK[$p - 1] =~ "0") ) {
            push @data1_seq, $$data1[$p];
        } elsif ( ($$STACK[$p] =~ "1") && ($$STACK[$p - 1] =~ "1") ) {
            $lastEntry = (scalar @data1_seq) - 1;
            @data1_seq[$lastEntry] .= $data1[$p];
        } else {
            next;
        }
    }

    ## get the names into an array
    sub extractFastName {
        my ($data1, $STACK, $data1_name) = @_;
        use strict;
        use warnings;
    }
}

```

```
#####
##### PLEASE CHANGE THE PATHS FOR PROGRAMS LOCATIONS AND RESULTS LOCATION #####
#####
##### 1 ##### Path where all Perl (scripts) programs are located on your computer:
my $pathprog = "/Volumes/MachD2/DKFZ/PERL/";

##### 2 ##### HERE is the path where you want your results to be:
my $resultPath = "/Volumes/MachD2/DKFZ/";

#####
##### USERS area ENDS #####
#####
print "\n\n=====
print "\n\nRUNNING PROGRAM SET 2 (P2-7) AFTER SORTING, FOR BARCODE NUMBER : ".$input."\n";
print " " For experiments using primers E01-E32\n";
print "\n\n The program set works on the REVERSE-COMPLEMENT sequences!!\n";

my $p2 = "set2_p2mac.pl";
my $p3ee = "set2_p3mac.pl";
my $p4 = "set2_p4mac.pl";
my $p5a = "set2_p5mac_a_noCutoff.pl";
my $p5d = "set2_p5mac_d_28bpCutoff.pl";
my $p5after = "set2_p5mac_edit.pl";
my $p6ee = "set2_p6mac_before.pl";
my $p6afteree = "set2_p6mac_fasta.pl";
my $p9ee = "set2_p7mac.pl";

system ("perl $pathprog$p2 $input $pyroNum $resultPath");
system ("perl $pathprog$p3ee $input $pyroNum $resultPath");
system ("perl $pathprog$p4 $input $pyroNum $resultPath");
system ("perl $pathprog$p5a $input $pyroNum $resultPath");
system ("perl $pathprog$p5d $input $pyroNum $resultPath");
system ("perl $pathprog$p5after $input $pyroNum $resultPath");
system ("perl $pathprog$p6ee $input $pyroNum $resultPath");
system ("perl $pathprog$p6afteree $input $pyroNum $resultPath");
system ("perl $pathprog$p9ee $input $pyroNum $resultPath");

print "\n\nRUNNING ALL PROGRAMS (2-7) AFTER SORTING, FOR BARCODE NUMBER : ".$input."\n";
print ".....FINISHED\n";
print "\n\n===== \n\n";
exit;

##x Main Program ENDS HERE xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
#####

#####
##### About the script #####
# Created by Sasithorn Chotewutmontri, Jan 2009.
#
#####
##### MAIN PROGRAM BODY #####
#####
use strict;
use warnings;

my @sampleNo = @ARGV;
my $input = $sampleNo[0];
my $pyroNum = $sampleNo[1];
```

Program set 3:

File name: **PROGRAM_SET_3_mac.pl**

Source code:

```
#!/usr/bin/perl

## About the script #####
# Created by Sasithorn Chotewutmontri, Jan 2009.
#
#####
##### MAIN PROGRAM BODY #####
#####
use strict;
use warnings;

my @sampleNo = @ARGV;
my $input = $sampleNo[0];
my $pyroNum = $sampleNo[1];
```

```
my $maxcounter = scalar @sdata1;
for (my $p = 0; $p < $maxcounter; $p++) {
    if ($sdata1[$p] =~ m/^0$/) {
        push @sdata1_name, $sdata1[$p];
    } else { next; }
}

}

## sorted by barcodes, the sequences are separated into individual sample
sub scanSortBarcode2 {
    my ($data1_name, $data1_seq, $sampleXX_name, $sampleXX_seq, $barcode, $code) = @_;
    use strict;
    use warnings;

    my $maxcount = scalar @sdata1_name;
    for (my $p = 0; $p < $maxcount; $p++) {
        if ($sdata1_seq[$p] =~ m/^$barcode$/) {
            push @sampleXX_name, $data1_name[$p];
            push @sampleXX_seq, $sdata1_seq[$p];
        } else { next; }
    }

}

## write the FastA formatted sequences of each sample into files in D:/perlresult/
sub saveSortSeqToFile {
    my ($sampleXX_name, $sampleXX_seq, $filename, $path, $dirname) = @_;
    use strict;
    use warnings;

    my @nameAndSeq = ();
    system ("mkdir", "$path.$dirname");
    my $maxcount = scalar @sampleXX_name;
    for (my $p = 0; $p < $maxcount; $p++) {
        push @nameAndSeq, $sampleXX_name[$p];
        push @nameAndSeq, $sampleXX_seq[$p];
    }
    open (NAMESEQ, ">$path.$dirname/$filename");
    print NAMESEQ @nameAndSeq;
    close NAMESEQ;

    my $nameList = "nameList";
    open (NAME, ">$path.$dirname/sample.$dirname.nameList");
    for (my $n = 0; $n < $maxcount; $n++) {
        print NAME $sampleXX_name[$n];
    }
    close NAME;
}

}

##x Sub-routines END HERE xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
#####

#####
##### About the script #####
# Created by Sasithorn Chotewutmontri, Jan 2009.
#
#####
##### MAIN PROGRAM BODY #####
#####
use strict;
use warnings;

my @sampleNo = @ARGV;
my $input = $sampleNo[0];
my $pyroNum = $sampleNo[1];
```

Program set 2:

File name: **PROGRAM_SET_2_mac.pl**

Source code:

```
#!/usr/bin/perl

## About the script #####
# Created by Sasithorn Chotewutmontri, Jan 2009.
#
#####
##### MAIN PROGRAM BODY #####
#####
use strict;
use warnings;

my @sampleNo = @ARGV;
my $input = $sampleNo[0];
my $pyroNum = $sampleNo[1];
```



```
my @allbarcodes = ("01","02","03","04","05","06","07","08","09","10","11","12",  
                  "13","14","15","16","17","18","19","20","21","22","23","24");  
  
print "\n\n=====";  
print "\n\nRUNNING PROGRAMS (10-11), FOR ALL barcodes:\n";  
print "      This program is required to run ONLY ONCE!\n";  
print "\n\n    It will perform for ALL BARCODES, even one barcode is given as input\n";  
  
my $pBefore22ee = "set4_p10mac.pl";  
my $p22ee = "set4_p11mac.pl";  
system ("perl $pathprog$pBefore22ee >input $pyrnum $resultPath");  
  
for (my $p = 0; $p < (scalar @allbarcodes); $p++) {  
    system ("perl $pathprog$p22ee >$allbarcodes[$p] $pyrnum $resultPath");  
}  
  
print "\n\nRUNNING ALL PROGRAMS (10-11), FOR ALL BARCODES\n";  
print "  
print "  
print "\n\n===== FINISHED\n";  
exit;  
  
##x Main Program ENDS HERE xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx  
#####
```

Sub-programs for program set 2-4:

Program set 2 consists of nine sub-programs whose names start with prefix “set2_”. Program set 3 consists of two sub-programs with prefix “set3_”. And program set 4 consists of the last two sub-programs with prefix “set4_”.

File name: **set2_p2mac.pl**

Source code:

```
#!/usr/bin/perl

#### About the script #####
# Created by Sasithorn Choteuwitmontri, Jan 2009.

#####
## MAIN PROGRAM BODY #####
#####

use strict;
use warnings;

my $sampleNo = @ARGV;

my $blast = "blast-2.2.17";
my $blastall = $blast."bin/blastall";
my $dbHPV = $blast."db/hpv16";

my $prefix_path= $sampleNo[2];
my $pyroNum = int($sampleNo[1]);

my $path = $prefix_path."_Pyro".$pyroNum."_result";
my $inputFile = $path.$sampleNo[0]."/sample".$sampleNo[0].".txt";
# INPUT FASTA
# This input contains sorted sequences acc. to barcode,
# They will be reverse-complemented before being blasted to hpv16

my $pathSave = $path.$sampleNo[0]."/blastHPV16";
my $infoFile = "numberOfsequences.txt";
my $wordKdIr = $path.$sampleNo[0].".";
my $src = $wordKdIr."sample".$sampleNo[0]."_RC.txt";

print "\n\nPROGRAM 2 (BLAST HPV16) STARTS.....\n\n";

system ("mkdir", $pathSave);

my @data2 = ();
```

```
#####  
##### PLEASE CHANGE THE PATHS FOR PROGRAM LOCATIONS AND RESULTS LOCATION #####  
##### USERS can change paths here #####  
##### Path where all Perl (scripts) programs are located on your computer:  
  
my $spathprog = "/Volumes/MachD02/DKFZ/PERL/";  
  
##### 1 #####  
  
##### HERE is the path where you want your results to be:  
  
my $resultPath = "/Volumes/MachD02/DKFZ/";  
  
##### Users area ENDS #####  
#####  
  
print "\n\n===== \\\n";  
print "python RUNNING PROGRAMS (8-9), FOR BARCODE NUMBER : ".$input."\n";  
print "      For experiments using primers E01-E32\n";  
print "      The program set works on the REVERSE-COMPLEMENT sequences!!\n";  
  
my $splice = "set3_p8mac.pl";  
my $splicee = "set3_p9mac.pl";  
  
system ("perl $pathprog$splice $input $pyrNum $resultPath");  
system ("perl $pathprog$splicee $input $pyrNum $resultPath");  
  
print "python RUNNING ALL PROGRAMS (8-9), FOR BARCODE NUMBER : ".$input."\n";  
print "python .....FINISHED!\n";  
print "\n\n=====\n";  
exit;
```

Program set 4:

File name: **PROGRAM SET 4 mac.pl**

Source code:

```
#!/usr/bin/perl #####  
##### About the script #####  
# # Created by Sasithorn Chotewutmontri, Jan 2009.  
# #####  
  
use strict;  
use warnings;  
  
my @sampleNo = @ARGV;  
my $inpnt = $sampleNo[0];  
my $pyroNum = $sampleNo[1];  
  
##### USERS can change paths here #####  
##### PLEASE CHANGE THE PATHS FOR PROGRAM LOCATIONS AND RESULTS LOCATION #####  
##### Path where all Perl (scripts) programs are located on your computer:  
my $pathprog = "/Volumes/MachO2/DKFZ/PERL/";  
  
##### 2 ##### HERE is the path where you want your results to be:  
my $resultPath = "/Volumes/MachO2/DKFZ/";
```

[illegible]

[illegible]

```
# perfect' & complete match (CASE A1)
if ( ($bl_6 = $plus) && (int($bl_4) == 00)
&& (int($bl_3) >= int($bl_2)) ) {

    if ($bl_1 =~ "E17") {
        sel+++, print E1 Query, "#E17#complete\n";
        print COMBI Query, "#E17#complete\n";
    } elseif ($bl_1 =~ "E02") {
        sel2++, print E2 Query, "#E02#complete\n";
        print COMBI Query, "#E02#complete\n";
    } elseif ($bl_1 =~ "E03") {
        sel3++, print E3 Query, "#E03#complete\n";
        print COMBI Query, "#E03#complete\n";
    } elseif ($bl_1 =~ "E04") {
        sel4++, print E4 Query, "#E04#complete\n";
        print COMBI Query, "#E04#complete\n";
    } elseif ($bl_1 =~ "E05") {
        sel5++, print E5 Query, "#E05#complete\n";
        print COMBI Query, "#E05#complete\n";
    } elseif ($bl_1 =~ "E06") {
        sel6++, print E6 Query, "#E06#complete\n";
        print COMBI Query, "#E06#complete\n";
    } elseif ($bl_1 =~ "E07") {
        sel7++, print E7 Query, "#E07#complete\n";
        print COMBI Query, "#E07#complete\n";
    } elseif ($bl_1 =~ "E08") {
        sel8++, print E8 Query, "#E08#complete\n";
        print COMBI Query, "#E08#complete\n";
    } elseif ($bl_1 =~ "E09") {
        sel9++, print L1 Query, "#E09#complete\n";
        print COMBI Query, "#E09#complete\n";
    } elseif ($bl_1 =~ "E21") {
        sel12++, print L2 Query, "#E21#complete\n";
        print COMBI Query, "#E21#complete\n";
    } elseif ($bl_1 =~ "E11") {
        sel13++, print L3 Query, "#E11#complete\n";
        print COMBI Query, "#E11#complete\n";
    } elseif ($bl_1 =~ "E12") {
        sel14++, print L4 Query, "#E12#complete\n";
        print COMBI Query, "#E12#complete\n";
    } elseif ($bl_1 =~ "E13") {
        sel15++, print L5 Query, "#E13#complete\n";
        print COMBI Query, "#E13#complete\n";
    } elseif ($bl_1 =~ "E14") {
        sel16++, print L6 Query, "#E14#complete\n";
        print COMBI Query, "#E14#complete\n";
    } elseif ($bl_1 =~ "E15") {
        sel17++, print L7 Query, "#E15#complete\n";
        print COMBI Query, "#E15#complete\n";
    } elseif ($bl_1 =~ "E16") {
        sel18++, print L8 Query, "#E16#complete\n";
        print COMBI Query, "#E16#complete\n";
    } elseif ($bl_1 =~ "E18") {
        sel19++, print H1 Query, "#E18#complete\n";
        print COMBI Query, "#E18#complete\n";
    } elseif ($bl_1 =~ "E19") {
        sel20++, print H2 Query, "#E19#complete\n";
        print COMBI Query, "#E19#complete\n";
    } elseif ($bl_1 =~ "E20") {
        sel21++, print H3 Query, "#E20#complete\n";
        print COMBI Query, "#E20#complete\n";
    } elseif ($bl_1 =~ "E22") {
        sel22++, print H4 Query, "#E22#complete\n";
        print COMBI Query, "#E22#complete\n";
    } elseif ($bl_1 =~ "E23") {
        sel23++, print H5 Query, "#E23#complete\n";
        print COMBI Query, "#E23#complete\n";
    } elseif ($bl_1 =~ "E24") {
        sel24++, print H6 Query, "#E24#complete\n";
        print COMBI Query, "#E24#complete\n";
    } elseif ($bl_1 =~ "E25") {
        sel25++, print H7 Query, "#E25#complete\n";
        print COMBI Query, "#E25#complete\n";
    } elseif ($bl_1 =~ "E26") {
        sel26++, print H8 Query, "#E26#complete\n";
        print COMBI Query, "#E26#complete\n";
    } elseif ($bl_1 =~ "E27") {

```

```

    $ln2++; print H2 $query."#E19#Incomplete_tilHPV\n";
    print COMBI $query."#E19#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E20") {
    $ln3++; print H3 $query."#E20#Incomplete_tilHPV\n";
    print COMBI $query."#E20#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E22") {
    $ln4++; print H4 $query."#E22#Incomplete_tilHPV\n";
    print COMBI $query."#E22#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E23") {
    $ln5++; print H5 $query."#E23#Incomplete_tilHPV\n";
    print COMBI $query."#E23#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E24") {
    $ln6++; print H6 $query."#E24#Incomplete_tilHPV\n";
    print COMBI $query."#E24#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E25") {
    $ln7++; print H7 $query."#E25#Incomplete_tilHPV\n";
    print COMBI $query."#E25#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E26") {
    $ln8++; print H8 $query."#E26#Incomplete_tilHPV\n";
    print COMBI $query."#E26#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E27") {
    $lnk1++; print K1 $query."#E27#Incomplete_tilHPV\n";
    print COMBI $query."#E27#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E28") {
    $lnk2++; print K2 $query."#E28#Incomplete_tilHPV\n";
    print COMBI $query."#E28#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E29") {
    $lnk3++; print K3 $query."#E29#Incomplete_tilHPV\n";
    print COMBI $query."#E29#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E30") {
    $lnk4++; print K4 $query."#E30#Incomplete_tilHPV\n";
    print COMBI $query."#E30#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E31") {
    $lnk5++; print K5 $query."#E31#Incomplete_tilHPV\n";
    print COMBI $query."#E31#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E32") {
    $lnk6++; print K6 $query."#E32#Incomplete_tilHPV\n";
    print COMBI $query."#E32#Incomplete_tilHPV\n";
  } else { next; }

# 'imperfect' & incomplete match (CASES B1/2, C1/2)
} elsif ( ($b1_6 =~ $plus) && (int($b1_4) >= 90)
&& (int($b1_3) < int($b1_2))
&& (int($b1_7) == 1)
&& (int($b1_3)+int($b1_7)) > 20
&& (int($b1_2)-(int($b1_3)+int($b1_7))) < 5 ) {
  if ($b1_1 =~ "E17") {
    $ln1++; print E1 $query."#E17#Incomplete_middle\n";
    print COMBI $query."#E17#Incomplete_middle\n";
  } elsif ($b1_1 =~ "E02") {
    $ln2++; print E2 $query."#E02#Incomplete_middle\n";
    print COMBI $query."#E02#Incomplete_middle\n";
  } elsif ($b1_1 =~ "E03") {
    $ln3++; print E3 $query."#E03#Incomplete_middle\n";
    print COMBI $query."#E03#Incomplete_middle\n";
  } elsif ($b1_1 =~ "E04") {
    $ln4++; print E4 $query."#E04#Incomplete_middle\n";
    print COMBI $query."#E04#Incomplete_middle\n";
  } elsif ($b1_1 =~ "E05") {
    $ln5++; print E5 $query."#E05#Incomplete_middle\n";
    print COMBI $query."#E05#Incomplete_middle\n";
  } elsif ($b1_1 =~ "E06") {
    $ln6++; print E6 $query."#E06#Incomplete_middle\n";
    print COMBI $query."#E06#Incomplete_middle\n";
  } elsif ($b1_1 =~ "E07") {
    $ln7++; print E7 $query."#E07#Incomplete_middle\n";
    print COMBI $query."#E07#Incomplete_middle\n";
  } elsif ($b1_1 =~ "E08") {
    $ln8++; print E8 $query."#E08#Incomplete_middle\n";
    print COMBI $query."#E08#Incomplete_middle\n";
  } elsif ($b1_1 =~ "E09") {
    $lnl1++; print L1 $query."#E09#Incomplete_middle\n";
    print COMBI $query."#E09#Incomplete_middle\n";
  } elsif ($b1_1 =~ "E21") {
    $lnl2++; print L2 $query."#E21#Incomplete_middle\n";
    print COMBI $query."#E21#Incomplete_middle\n";
  }
}

$ln1++; print K1 $query."#E27#complete\n";
print COMBI $query."#E27#complete\n";
} elsif ($b1_1 =~ "E28") {
  $ln2++; print K2 $query."#E28#complete\n";
  print COMBI $query."#E28#complete\n";
} elsif ($b1_1 =~ "E29") {
  $ln3++; print K3 $query."#E29#complete\n";
  print COMBI $query."#E29#complete\n";
} elsif ($b1_1 =~ "E30") {
  $ln4++; print K4 $query."#E30#complete\n";
  print COMBI $query."#E30#complete\n";
} elsif ($b1_1 =~ "E31") {
  $ln5++; print K5 $query."#E31#complete\n";
  print COMBI $query."#E31#complete\n";
} elsif ($b1_1 =~ "E32") {
  $ln6++; print K6 $query."#E32#complete\n";
  print COMBI $query."#E32#complete\n";
} else { next; }

# 'imperfect' & complete match (CASE A2)
} elsif ( ($b1_6 =~ $plus) && (int($b1_4)
&& (int($b1_3) < int($b1_2))
&& (int($b1_7) == 1)
&& ((int($b1_3)+int($b1_7)) > 20)
&& ((int($b1_2)-(int($b1_3)+int($b1_7))) < 5) ) {
  if ($b1_1 =~ "E17") {
    $ln1++; print E1 $query."#E17#Incomplete_tilHPV\n";
    print COMBI $query."#E17#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E02") {
    $ln2++; print E2 $query."#E02#Incomplete_tilHPV\n";
    print COMBI $query."#E02#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E03") {
    $ln3++; print E3 $query."#E03#Incomplete_tilHPV\n";
    print COMBI $query."#E03#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E04") {
    $ln4++; print E4 $query."#E04#Incomplete_tilHPV\n";
    print COMBI $query."#E04#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E05") {
    $ln5++; print E5 $query."#E05#Incomplete_tilHPV\n";
    print COMBI $query."#E05#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E06") {
    $ln6++; print E6 $query."#E06#Incomplete_tilHPV\n";
    print COMBI $query."#E06#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E07") {
    $ln7++; print E7 $query."#E07#Incomplete_tilHPV\n";
    print COMBI $query."#E07#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E08") {
    $ln8++; print E8 $query."#E08#Incomplete_tilHPV\n";
    print COMBI $query."#E08#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E09") {
    $lnl1++; print L1 $query."#E09#Incomplete_tilHPV\n";
    print COMBI $query."#E09#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E21") {
    $lnl2++; print L2 $query."#E21#Incomplete_tilHPV\n";
    print COMBI $query."#E21#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E13") {
    $ln13++; print L3 $query."#E13#Incomplete_tilHPV\n";
    print COMBI $query."#E13#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E12") {
    $ln14++; print L4 $query."#E12#Incomplete_tilHPV\n";
    print COMBI $query."#E12#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E13") {
    $ln15++; print L5 $query."#E13#Incomplete_tilHPV\n";
    print COMBI $query."#E13#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E14") {
    $ln16++; print L6 $query."#E14#Incomplete_tilHPV\n";
    print COMBI $query."#E14#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E15") {
    $ln17++; print L7 $query."#E15#Incomplete_tilHPV\n";
    print COMBI $query."#E15#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E16") {
    $ln18++; print L8 $query."#E16#Incomplete_tilHPV\n";
    print COMBI $query."#E16#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E18") {
    $lnl1++; print L1 $query."#E18#Incomplete_tilHPV\n";
    print COMBI $query."#E18#Incomplete_tilHPV\n";
  } elsif ($b1_1 =~ "E19") {
    $lnl2++; print L2 $query."#E19#Incomplete_tilHPV\n";
    print COMBI $query."#E19#Incomplete_tilHPV\n";
  }
}

```

```

    } elseif ($b1_1 =~ "E11") {
        $n66++; print E6 query,"#E06#00T\n"; print COMBI query,"#E06#00T\n";
    } elseif ($b1_1 =~ "E07") {
        $n67++; print E7 query,"#E07#00T\n"; print COMBI query,"#E07#00T\n";
    } elseif ($b1_1 =~ "E08") {
        $n68++; print E8 query,"#E08#00T\n"; print COMBI query,"#E08#00T\n";
    } elseif ($b1_1 =~ "E09") {
        $n69++; print E9 query,"#E09#00T\n"; print COMBI query,"#E09#00T\n";
    } elseif ($b1_1 =~ "E10") {
        $n70++; print E10 query,"#E10#00T\n"; print COMBI query,"#E10#00T\n";
    } elseif ($b1_1 =~ "E11") {
        $n71++; print E11 query,"#E11#00T\n"; print COMBI query,"#E11#00T\n";
    } elseif ($b1_1 =~ "E12") {
        $n72++; print E12 query,"#E12#00T\n"; print COMBI query,"#E12#00T\n";
    } elseif ($b1_1 =~ "E13") {
        $n73++; print E13 query,"#E13#00T\n"; print COMBI query,"#E13#00T\n";
    } elseif ($b1_1 =~ "E14") {
        $n74++; print E14 query,"#E14#00T\n"; print COMBI query,"#E14#00T\n";
    } elseif ($b1_1 =~ "E15") {
        $n75++; print E15 query,"#E15#00T\n"; print COMBI query,"#E15#00T\n";
    } elseif ($b1_1 =~ "E16") {
        $n76++; print E16 query,"#E16#00T\n"; print COMBI query,"#E16#00T\n";
    } elseif ($b1_1 =~ "E17") {
        $n77++; print E17 query,"#E17#00T\n"; print COMBI query,"#E17#00T\n";
    } elseif ($b1_1 =~ "E18") {
        $n78++; print E18 query,"#E18#00T\n"; print COMBI query,"#E18#00T\n";
    } elseif ($b1_1 =~ "E19") {
        $n79++; print E19 query,"#E19#00T\n"; print COMBI query,"#E19#00T\n";
    } elseif ($b1_1 =~ "E20") {
        $n80++; print E20 query,"#E20#00T\n"; print COMBI query,"#E20#00T\n";
    } elseif ($b1_1 =~ "E21") {
        $n81++; print E21 query,"#E21#00T\n"; print COMBI query,"#E21#00T\n";
    } elseif ($b1_1 =~ "E22") {
        $n82++; print E22 query,"#E22#00T\n"; print COMBI query,"#E22#00T\n";
    } elseif ($b1_1 =~ "E23") {
        $n83++; print E23 query,"#E23#00T\n"; print COMBI query,"#E23#00T\n";
    } elseif ($b1_1 =~ "E24") {
        $n84++; print E24 query,"#E24#00T\n"; print COMBI query,"#E24#00T\n";
    } elseif ($b1_1 =~ "E25") {
        $n85++; print E25 query,"#E25#00T\n"; print COMBI query,"#E25#00T\n";
    } elseif ($b1_1 =~ "E26") {
        $n86++; print E26 query,"#E26#00T\n"; print COMBI query,"#E26#00T\n";
    } elseif ($b1_1 =~ "E27") {
        $n87++; print E27 query,"#E27#00T\n"; print COMBI query,"#E27#00T\n";
    } elseif ($b1_1 =~ "E28") {
        $n88++; print E28 query,"#E28#00T\n"; print COMBI query,"#E28#00T\n";
    } elseif ($b1_1 =~ "E29") {
        $n89++; print E29 query,"#E29#00T\n"; print COMBI query,"#E29#00T\n";
    } elseif ($b1_1 =~ "E30") {
        $n90++; print E30 query,"#E30#00T\n"; print COMBI query,"#E30#00T\n";
    } elseif ($b1_1 =~ "E31") {
        $n91++; print E31 query,"#E31#00T\n"; print COMBI query,"#E31#00T\n";
    } elseif ($b1_1 =~ "E32") {
        $n92++; print E32 query,"#E32#00T\n"; print COMBI query,"#E32#00T\n";
    } else {
        $n93++; print E33 query,"#E33#00T\n"; print COMBI query,"#E33#00T\n";
    }
}

$hitOutranged++;
if ($b1_1 =~ "E17") {
    $n94++; print E1 query,"#E17#00T_ranged\n";
    print COMBI query,"#E17#00T_ranged\n";
} elseif ($b1_1 =~ "E02") {
    $n95++; print E2 query,"#E02#00T_ranged\n";
    print COMBI query,"#E02#00T_ranged\n";
} elseif ($b1_1 =~ "E03") {
    $n96++; print E3 query,"#E03#00T_ranged\n";
    print COMBI query,"#E03#00T_ranged\n";
} elseif ($b1_1 =~ "E04") {
    $n97++; print E4 query,"#E04#00T_ranged\n";
    print COMBI query,"#E04#00T_ranged\n";
} elseif ($b1_1 =~ "E05") {
    $n98++; print E5 query,"#E05#00T_ranged\n";
    print COMBI query,"#E05#00T_ranged\n";
} elseif ($b1_1 =~ "E06") {
    $n99++; print E6 query,"#E06#00T_ranged\n";
    print COMBI query,"#E06#00T_ranged\n";
} elseif ($b1_1 =~ "E07") {
    $n100++; print E7 query,"#E07#00T_ranged\n";
    print COMBI query,"#E07#00T_ranged\n";
} elseif ($b1_1 =~ "E08") {
    $n101++; print E8 query,"#E08#00T_ranged\n";
    print COMBI query,"#E08#00T_ranged\n";
} else {
    $n102++; print E9 query,"#E09#00T_ranged\n";
    print COMBI query,"#E09#00T_ranged\n";
}
}

} elseif ($b1_1 =~ "E11") {
    $n103++; print L3 query,"#E11#Incomplete_middl\n";
    print COMBI query,"#E11#Incomplete_middl\n";
} elseif ($b1_1 =~ "E12") {
    $n104++; print L4 query,"#E12#Incomplete_middl\n";
    print COMBI query,"#E12#Incomplete_middl\n";
} elseif ($b1_1 =~ "E13") {
    $n105++; print L5 query,"#E13#Incomplete_middl\n";
    print COMBI query,"#E13#Incomplete_middl\n";
} elseif ($b1_1 =~ "E14") {
    $n106++; print L6 query,"#E14#Incomplete_middl\n";
    print COMBI query,"#E14#Incomplete_middl\n";
} elseif ($b1_1 =~ "E15") {
    $n107++; print L7 query,"#E15#Incomplete_middl\n";
    print COMBI query,"#E15#Incomplete_middl\n";
} elseif ($b1_1 =~ "E16") {
    $n108++; print L8 query,"#E16#Incomplete_middl\n";
    print COMBI query,"#E16#Incomplete_middl\n";
} elseif ($b1_1 =~ "E17") {
    $n109++; print H1 query,"#E17#Incomplete_middl\n";
    print COMBI query,"#E17#Incomplete_middl\n";
} elseif ($b1_1 =~ "E18") {
    $n110++; print H2 query,"#E18#Incomplete_middl\n";
    print COMBI query,"#E18#Incomplete_middl\n";
} elseif ($b1_1 =~ "E19") {
    $n111++; print H3 query,"#E19#Incomplete_middl\n";
    print COMBI query,"#E19#Incomplete_middl\n";
} elseif ($b1_1 =~ "E20") {
    $n112++; print H4 query,"#E20#Incomplete_middl\n";
    print COMBI query,"#E20#Incomplete_middl\n";
} elseif ($b1_1 =~ "E21") {
    $n113++; print H5 query,"#E21#Incomplete_middl\n";
    print COMBI query,"#E21#Incomplete_middl\n";
} elseif ($b1_1 =~ "E22") {
    $n114++; print H6 query,"#E22#Incomplete_middl\n";
    print COMBI query,"#E22#Incomplete_middl\n";
} elseif ($b1_1 =~ "E23") {
    $n115++; print H7 query,"#E23#Incomplete_middl\n";
    print COMBI query,"#E23#Incomplete_middl\n";
} elseif ($b1_1 =~ "E24") {
    $n116++; print H8 query,"#E24#Incomplete_middl\n";
    print COMBI query,"#E24#Incomplete_middl\n";
} elseif ($b1_1 =~ "E25") {
    $n117++; print H9 query,"#E25#Incomplete_middl\n";
    print COMBI query,"#E25#Incomplete_middl\n";
} elseif ($b1_1 =~ "E26") {
    $n118++; print H10 query,"#E26#Incomplete_middl\n";
    print COMBI query,"#E26#Incomplete_middl\n";
} elseif ($b1_1 =~ "E27") {
    $n119++; print H11 query,"#E27#Incomplete_middl\n";
    print COMBI query,"#E27#Incomplete_middl\n";
} elseif ($b1_1 =~ "E28") {
    $n120++; print H12 query,"#E28#Incomplete_middl\n";
    print COMBI query,"#E28#Incomplete_middl\n";
} elseif ($b1_1 =~ "E29") {
    $n121++; print H13 query,"#E29#Incomplete_middl\n";
    print COMBI query,"#E29#Incomplete_middl\n";
} elseif ($b1_1 =~ "E30") {
    $n122++; print H14 query,"#E30#Incomplete_middl\n";
    print COMBI query,"#E30#Incomplete_middl\n";
} elseif ($b1_1 =~ "E31") {
    $n123++; print H15 query,"#E31#Incomplete_middl\n";
    print COMBI query,"#E31#Incomplete_middl\n";
} elseif ($b1_1 =~ "E32") {
    $n124++; print H16 query,"#E32#Incomplete_middl\n";
    print COMBI query,"#E32#Incomplete_middl\n";
} elseif ($b1_1 =~ "E33") {
    $n125++; print H17 query,"#E33#Incomplete_middl\n";
    print COMBI query,"#E33#Incomplete_middl\n";
} else {
    $n126++; print H18 query,"#E34#Incomplete_middl\n";
    print COMBI query,"#E34#Incomplete_middl\n";
}
}

} elseif ( (( $b1_6 =~ $p1us ) && (int($b1_4) < 90) ) |
($b1_6 =~ $m1us) ) {
    $hitOutranged++;
    if ($b1_1 =~ "E17") {
        $n127++; print E1 query,"#E17#00T\n"; print COMBI query,"#E17#00T\n";
    } elseif ($b1_1 =~ "E02") {
        $n128++; print E2 query,"#E02#00T\n"; print COMBI query,"#E02#00T\n";
    } elseif ($b1_1 =~ "E03") {
        $n129++; print E3 query,"#E03#00T\n"; print COMBI query,"#E03#00T\n";
    } elseif ($b1_1 =~ "E04") {
        $n130++; print E4 query,"#E04#00T\n"; print COMBI query,"#E04#00T\n";
    } elseif ($b1_1 =~ "E05") {
        $n131++; print E5 query,"#E05#00T\n"; print COMBI query,"#E05#00T\n";
    } elseif ($b1_1 =~ "E06") {
        $n132++; print E6 query,"#E06#00T\n"; print COMBI query,"#E06#00T\n";
    }
}
}

```

```
} elseif ($d1_1 == "E09") {  
    $n11+=; print L1 Query,"#E09$OUT_ranged\n";  
    print COMB1 query,"#E09$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E21") {  
    $n12+=; print L2 Query,"#E21$OUT_ranged\n";  
    print COMB1 query,"#E21$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E11") {  
    $n13+=; print L3 Query,"#E11$OUT_ranged\n";  
    print COMB1 query,"#E11$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E12") {  
    $n14+=; print L4 Query,"#E12$OUT_ranged\n";  
    print COMB1 query,"#E12$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E13") {  
    $n15+=; print L5 Query,"#E13$OUT_ranged\n";  
    print COMB1 query,"#E13$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E14") {  
    $n16+=; print L6 Query,"#E14$OUT_ranged\n";  
    print COMB1 query,"#E14$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E15") {  
    $n17+=; print L7 Query,"#E15$OUT_ranged\n";  
    print COMB1 query,"#E15$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E16") {  
    $n18+=; print L8 Query,"#E16$OUT_ranged\n";  
    print COMB1 query,"#E16$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E18") {  
    $n19+=; print H1 Query,"#E18$OUT_ranged\n";  
    print COMB1 query,"#E18$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E19") {  
    $n20+=; print H2 Query,"#E19$OUT_ranged\n";  
    print COMB1 query,"#E19$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E20") {  
    $n21+=; print H3 Query,"#E20$OUT_ranged\n";  
    print COMB1 query,"#E20$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E22") {  
    $n22+=; print H4 Query,"#E22$OUT_ranged\n";  
    print COMB1 query,"#E22$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E23") {  
    $n23+=; print H5 Query,"#E23$OUT_ranged\n";  
    print COMB1 query,"#E23$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E24") {  
    $n24+=; print H6 Query,"#E24$OUT_ranged\n";  
    print COMB1 query,"#E24$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E25") {  
    $n25+=; print H7 Query,"#E25$OUT_ranged\n";  
    print COMB1 query,"#E25$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E26") {  
    $n26+=; print H8 Query,"#E26$OUT_ranged\n";  
    print COMB1 query,"#E26$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E27") {  
    $n27+=; print K1 Query,"#E27$OUT_ranged\n";  
    print COMB1 query,"#E27$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E28") {  
    $n28+=; print K2 Query,"#E28$OUT_ranged\n";  
    print COMB1 query,"#E28$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E29") {  
    $n29+=; print K3 Query,"#E29$OUT_ranged\n";  
    print COMB1 query,"#E29$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E30") {  
    $n30+=; print K4 Query,"#E30$OUT_ranged\n";  
    print COMB1 query,"#E30$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E31") {  
    $n31+=; print K5 Query,"#E31$OUT_ranged\n";  
    print COMB1 query,"#E31$OUT_ranged\n";  
}  
} elseif ($d1_1 == "E32") {  
    $n32+=; print K6 Query,"#E32$OUT_ranged\n";  
    print COMB1 query,"#E32$OUT_ranged\n";  
}  
} else { next; }  
  
} else { next; }  
  
}  
  
unless ($open ">$out")) { print "\n\n\nCan not create \"$out.\n\n\n"; exit;  
}$txt1 = : complete vs Incomplete &\nOutraged \</code>
```


File name: **set2_p4mac.pl**

Source code:

```
#!/usr/bin/perl
## About the script #####
# Created by Sasithorn Chotewutmontri, Jan 2009.
#
#####
### MAIN PROGRAM BODY #####
#####
use strict;
use warnings;

my $sampleNo = @ARGV;
my $prefix_path= $sampleNo[2];
my $pyroNum = int($sampleNo[1]);

my $path = $prefix_path."Pyro".$pyroNum."_result/";

my $blastFilePath = $path.$sampleNo[0]."/blastHPV16/";
my $workDir = $path.$sampleNo[0]."/";

my $numberOfseqFile = $workDir."numberOfsequences.txt";
my $outfile = $blastFilePath."parsed_hpv_allreport";
my $outstat = $blastFilePath."parsed_hpv_statistics";

print "\n\nPROGRAM 4 (PARSING all) STARTS.....";

my @totSeq = getFiledData ($numberOfseqfile);
my $count = int($totSeq[0]);

getParseDouble2 ($blastFilePath, $count, $outfile, $outstat);

print "FINISHED\n\n";

## Main Program ENDS HERE
#####
### SUBROUTINES #####
#####
sub getFiledData {
    my ($filename) = @_;
    use strict;
    use warnings;

    my @filedata = ();
    unless ( open(FILE_DATA, $filename) ) {
        print STDERR "\n\nThe program can't open the files \"$filename\"!\n\n";
        print "Please re-check input file name and its location\n";
        print "The correct command should be :\n\n";
        print "\t\tperl PROGRAM(locationname) INPUTFastAFILE(locationname)\n\n";
        print "The command should be given under directory";
        print " c:/Perl/bin> in case of Dos Terminal\n\n";
        exit;
    }
    @filedata = <FILE_DATA>;
    close FILE_DATA;
    return @filedata;
}

# for Part 1
sub getParseDouble2 {
    my ($blastpath, $count, $out, $stat) = @_;
    use strict;
    use warnings;
    unless (open (OUT, ">$out")) {
```

```
print "\n\nCan not create ".$out."\n\n\n";
    exit;
}

my ($nohitcount, $shitcount) = (0, 0);
for (my $p = 0; $p < $count; $p++) {
    my @blast = getFiledData ($blastpath."b_blasted_".$p);
    $blast[0] =~ m/^(.....\(((0-9)*)(.*/);
    # Query Length
    my $QL = $1;

    if ($blast[16] =~ m/^(.....No.hits.found/) {
        $nohitcount++;
        print OUT $p."#".$QL."#noHitHPV#noHPV\n";
    } else {
        $shitcount++;
        print OUT $p."#".$QL."#hitHPV";
        for (my $m = 0; $m < (scalar @blast); $m++) {
            if ($blast[$m] =~ m/^( Score(.*))/) {
                my @tmpIden = split (/^s+/, $blast[$m+1]);
                #scoreBP
                my $scoreBox = $tmpIden[3];
                $scoreBox =~ m/^(0-9)*\((0-9)*$/;
                my $scoreBP = $2;
                #percent match
                my $percentBox = $tmpIden[4];
                $percentBox =~ m/^(0-9)*$/;
                my $percent = $1;
                #strand
                $blast[$m+2] =~ m/^(Strand=Plus...(.*)$/;
                my $strand = $1;
                # first-Match on Query
                $blast[$m+5] =~ m/^(Query..(0-9)*)(.*)\((0-9)*)(.*)$/;
                my $firstMatchQue = $1;
                # first-Match on HPV16 (can be last position in case MINUS strand)
                $blast[$m+7] =~ m/^(Sbjct..(0-9)*)(.*)\((0-9)*)(.*)$/;
                my $firstMatchHpv = $1;

                print OUT "#".$scoreBP."$percent.%".
                    $firstMatchQue."$strand."$firstMatchHpv;

            } else { next; }
        }
        print OUT "\n";
    }
}

} # close for loop
close OUT;

unless (open (STAT, ">$stat")) {print "\n\nCan not create ".$stat."\n\n"; exit;}

print STAT "Number of sequences WITH HIT to HPV16 (SIG & INSIGNificant) : ".$shitcount."\n";
print STAT "Number of sequences WITH NO HIT to HPV16 (by BLASTN) : ".$nohitcount."\n";
print STAT "Number of total sequences having this barcode : ".$count."\n";

close STAT;
}

#x Sub-routines END HERE
#####
#####
```

File name: **set2_p5mac_a_noCutoff.p1**

Source code:

```
#!/usr/bin/perl

## About the script #####
# Created by Sasithorn Chotewutmontri, Jan 2009.
#####

#####
##### MAIN PROGRAM BODY #####
#####

use strict;
use warnings;

my $sampleNo = @ARGV;
my $prefixPath = $sampleNo[2];
my $pyroNum = int($sampleNo[1]);

my $path = $prefixPath."Pyro".$pyroNum."_result/";
my $workingDir = $path.$sampleNo[0]."/";
my $blastFilePath = $workingDir."blastHPV16/";

# For different cutoff conditions, change these variables:
# $newDir, $prefix

my $newDir = $workingDir."noCutoff/";
my $prefix = "noCutoff_";

my $selectRep = $prefix."select_report";
my $grp1 = $prefix."select_group1";
my $grp2 = $prefix."select_group2";
my $grp3 = $prefix."select_group3";
my $grp4 = $prefix."select_group4";
my $selectRepub = $prefix."select_report_doubleHPV";

my $reportallpos = $prefix."allreport_fullpieceposition";
my $scf2 = 94; # CUTOFF PERCENT MATCH FOR 2nd HIT!
my $scf3 = 15; # CUTOFF SCOREBP FOR 2nd HIT!

##### Start of PART 1 : SELECT into 4 groups #####
my $parsedFile = $blastFilePath."parsed_hpv_allreport";

my @allParseRep = ();

print "\n\nPROGRAM 5A (SELECT NO CUTOFF) STARTS.....\n\n";
print "\n\n This program arrange each sequence in this sample group into \n\n";
print "\n\n 4 'significants/non-significant HPV16 bit' groups (because no CF is used)\n\n";
print "\n\n according to its 'non-HPV16' part\n\n";
print "\n\n G1 >= the non-HPV is >= 150 nt\n\n";
print "\n\n G2 >= the non-HPV is >= 100 AND < 150 nt\n\n";
print "\n\n G3 >= the non-HPV is >= 50 AND < 100 nt\n\n";
print "\n\n G4 >= the non-HPV is >= 0 AND < 50 nt\n\n";

# read parsed data for selection
@allParseRep = getfileData($parsedFile);

# make Dir for selected result of 'no cut-off'-program
System ("mkdir", $newDir);

unless (open (SELECT, ">$newDir/selectRep")) {print "\n\nCan not create ".$newDir.$selectRep."\n\n"; exit;}
unless (open (GRP1, ">$newDir/grp1")) {print "\n\nCan not create ".$newDir.$grp1."\n\n"; exit;}
unless (open (GRP2, ">$newDir/grp2")) {print "\n\nCan not create ".$newDir.$grp2."\n\n"; exit;}
unless (open (GRP3, ">$newDir/grp3")) {print "\n\nCan not create ".$newDir.$grp3."\n\n"; exit;}
unless (open (GRP4, ">$newDir/grp4")) {print "\n\nCan not create ".$newDir.$grp4."\n\n"; exit;}
unless (open (DUB, ">$newDir/selectRepub")) {print "\n\nCan not create ".$newDir.$selectRepub."\n\n"; exit;}
unless (open (RALL, ">$newDir/reportallpos")) {print "\n\nCan not create ".$newDir.$reportallpos."\n\n"; exit;}

for ( my $p = 0 ; $p < (scalar @allParseRep) ; $p++) {
    my $tmp_line = split (/./, $allParseRep[$p]);
    my $num = $tmp_line[0];
    my $SQL = $tmp_line[1];

    # query's seqNo
    # query's length
}
```

```
my $tmp_index = $tmp_line[2];
my $blockCount = scalar (@tmp_line); # hit/nohit indice

my ($tmp_scoreBP, $tmp_1stMatch, $tmp_strand, $tmp_1stMatchS, $suse);
my ($scoreBP2, $firstM2, $strand2, $firstMS2);
my ($tmp_lastpos); # the position (non-hpv) next to the last hpv match position
my ($scoreBP1, $firstM1, $strand1, $firstMS1);
my ($series1, $series2);

# This first IF-loop defines the 'virtual matching area of HPV16'
# The value of the 'first position' of HPV-matching area will be specified
# as well as other values, e.g. 'virtual scoreBP' in case of two continuous hits

# ONE-hit only
if ((($tmp_index =~ "hitHPV") && ($blockCount == 4)) { # for ONE-HPV-HIT query
    my @tmpBox = split ( /%/, $tmp_line[3]); # scoreBP
    $tmp_scoreBP = $tmpBox[0]; # 1st-Match-position-on-QUERY-strand
    $tmp_1stMatch = $tmpBox[1]; # strand of the matched HPV
    $tmp_strand = $tmpBox[3]; # 1st-Match-position-on-Subj-strand
    $suse = "hitx1_userINST";

    $$scoreBP1 = $tmpBox[0];

    $firstM1 = $tmpBox[2];
    $strand1 = $tmpBox[3];
    $firstMS1 = $tmpBox[4];
    $series1 = "use_Alone";

    $$scoreBP2 = "none";

    $firstM2 = "none";
    $strand2 = "none";
    $firstMS2 = "none";
    $series2 = "none";

    } elsif (($tmp_index =~ "hitHPV") && ($blockCount > 4)) { # for MORE-HPV-HITS query
        my @tmpBox1 = split ( /%/, $tmp_line[3]); # scoreBP of FIRST HIT
        my $scoreBP1 = $tmpBox1[0]; # percent match of FIRST HIT
        my $percent1 = $tmpBox1[1]; # 1st-Match-position-on-QUERY-strand of FIRST HIT
        my $tmp1stMatch1 = $tmpBox1[2]; # strand of the matched HPV of FIRST HIT
        my $strand1 = $tmpBox1[3];
        my $tmp1stMatchS1 = $tmpBox1[4];

        my @tmpBox2 = split ( /%/, $tmp_line[4]); # scoreBP of SECOND HIT
        my $scoreBP2 = $tmpBox2[0]; # percent match of SECOND HIT
        my $percent2 = $tmpBox2[1]; # 1st-Match-position-on-QUERY-strand of SECOND HIT
        my $tmp2stMatch2 = $tmpBox2[2]; # strand of the matched HPV of SECOND HIT
        my $strand2 = $tmpBox2[3];
        my $tmp2stMatchS2 = $tmpBox2[4];

        # (1) : if ending-of-first-hit is EARLIER than ending-of-second-hit
        if ((int($tmp1stMatch1)+int($scoreBP1)) < (int($tmp2stMatch2)+int($scoreBP2))) {
            # (1.1): first & second hits begin at the same pos
            if ( int($tmp1stMatch1) == int($tmp2stMatch2) ) {
                $tmp_scoreBP = $scoreBP2;
                $tmp_strand = $strand2;
                $tmp_1stMatch = $tmp1stMatch2;
                $suse = "hitx2_useSECOND_1stHitAsSubset";
                $tmp_1stMatchS = $tmp1stMatchS2;

                $scoreBP1 = $scoreBP2;
                $firstM1 = $tmp1stMatch1;
                $strand1 = $tmp1stMatch1;
                $firstMS1 = int($tmp1stMatchS1);
                $series1 = "subset";

                $scoreBP2 = $scoreBP2;
                $firstM2 = $tmp1stMatch2;
                $strand2 = $tmp1stMatch2;
                $firstMS2 = int($tmp1stMatchS2);
                $series2 = "use_Alone";

            } # (1.2): first hit begins first
            } elsif ( int($tmp1stMatch1) < int($tmp1stMatch2) ) {
                # (1.2): first hit begins first
            }
        }
    }
}
```

```

$strand1 = $strand1;
$firstMS1 = int($t1stMatchS1);
$series1 = "use_Alonge";

$scoreBP2 = $scoreBP2;
$firstM2 = $t1stMatch2;
$strand2 = $strand2;
$firstMS2 = int($t1stMatchS2);
$series2 = "smallerThan_".$cf3;
}

# (1.3): second hit begins first
} else {
    $temp_scoreBP = $scoreBP2;
    $temp_1stMatch = $t1stMatch2;
    $temp_strand = $strand2;
    $use = "hitx2_useSECOND_1stHitasSubset";
    $temp_1stMatchS = $t1stMatchS2;

    $scoreBP1 = $scoreBP1;
    $firstM1 = $t1stMatch1;
    $strand1 = $strand1;
    $firstMS1 = int($t1stMatchS1);
    $series1 = "subset";

    $scoreBP2 = $scoreBP2;
    $firstM2 = $t1stMatch2;
    $strand2 = $strand2;
    $firstMS2 = int($t1stMatchS2);
    $series2 = "use_Alonge";
}

# (2) : if ending-of-first-hit is THE SAME as ending-of-second-hit
} elseif ((int($t1stMatch1)+int($scoreBP1))==(int($t1stMatch2)+int($scoreBP2))) {
    # (2.1): first & second hits begin at the same pos
    if ( int($t1stMatch1) == int($t1stMatch2) ) {
        $temp_scoreBP = $scoreBP1;
        $temp_1stMatch = $t1stMatch1;
        $temp_strand = $strand1;
        $use = "hitx2_useFIRST_2ndHitasSubset";
        $temp_1stMatchS = $t1stMatchS1;

        $scoreBP1 = $scoreBP1;
        $firstM1 = $t1stMatch1;
        $strand1 = $strand1;
        $firstMS1 = int($t1stMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $scoreBP2;
        $firstM2 = $t1stMatch2;
        $strand2 = $strand2;
        $firstMS2 = int($t1stMatchS2);
        $series2 = "subset";
    }

    # (2.2): first hit begins first
    } elseif ( int($t1stMatch1) < int($t1stMatch2) ) {
        $temp_scoreBP = $scoreBP1;
        $temp_1stMatch = $t1stMatch1;
        $temp_strand = $strand1;
        $use = "hitx2_useFIRST_2ndHitasSubset";
        $temp_1stMatchS = $t1stMatchS1;

        $scoreBP1 = $scoreBP1;
        $firstM1 = $t1stMatch1;
        $strand1 = $strand1;
        $firstMS1 = int($t1stMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $scoreBP2;
        $firstM2 = $t1stMatch2;
        $strand2 = $strand2;
        $firstMS2 = int($t1stMatchS2);
        $series2 = "subset";
    }
}

```

```

# If last match position of 1stHit is >= 1stMatch of 2.Hit (+2)
# (+2) is the allowance for the gap or if there is no gap then 2 nt distance
# between the two hits
if ( (int($t1stMatch1)+int($scoreBP1)) >= (int($t1stMatch2)+2) ) {
    # Control if the 2.Hit is long enough to be significantly counted
    if ((int($scoreBP2) >= $cf3 ) && (int($percent2) >= $cf2 )) {
        $use = "hitx2_useFIRSTInSECOND_continuous";
        # DoubleHPV-- independent of HPV16's 4 groups
        print DUB $Qnum."#".SQL."#hitHPV#". $use."\\n";
        $temp_scoreBP = int($scoreBP2)+(int($t1stMatch2)-int($t1stMatch1));
        $temp_1stMatch = $t1stMatch1;
        $temp_strand = $strand1;
        $scoreBP2 = $scoreBP2;
        $firstM2 = $t1stMatch2;
        $strand2 = $strand2;
        $temp_1stMatchS = $t1stMatchS1;

        $scoreBP1 = $scoreBP1;
        $firstM1 = $t1stMatch1;
        $strand1 = $strand1;
        $firstMS1 = int($t1stMatchS1);
        $series1 = "use_1st";

        $firstMS2 = int($t1stMatchS2);
        $series2 = "use_2nd";

        $use = "hitx2_useFIRST_tooSmall2ndHit";
        $temp_scoreBP = $scoreBP1;
        $temp_1stMatch = $t1stMatch1;
        $temp_strand = $strand1;
        $temp_1stMatchS = $t1stMatchS1;

        $scoreBP1 = $scoreBP1;
        $firstM1 = $t1stMatch1;
        $strand1 = $strand1;
        $firstMS1 = int($t1stMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $scoreBP2;
        $firstM2 = $t1stMatch2;
        $strand2 = $strand2;
        $firstMS2 = int($t1stMatchS2);
        $series2 = "smallerThan_".$cf3;
    }
}

} else {
    if ((int($scoreBP2) >= $cf3 ) && (int($percent2) >= $cf2 )) {
        $use = "hitx2_useFIRSTInSECOND_discontinuous";
        # DoubleHPV-- independent of HPV16's 4 groups
        print DUB $Qnum."#".SQL."#hitHPV#". $use."\\n";
        $temp_scoreBP = $scoreBP1; # DIFFERENT FROM "continuous"
        $temp_1stMatch = $t1stMatch1;
        $temp_strand = $strand1;
        $scoreBP2 = $scoreBP2;
        $firstM2 = $t1stMatch2;
        $strand2 = $strand2;
        $temp_1stMatchS = $t1stMatchS1;

        $scoreBP1 = $scoreBP1;
        $firstM1 = $t1stMatch1;
        $strand1 = $strand1;
        $firstMS1 = int($t1stMatchS1);
        $series1 = "use_1st";

        $firstMS2 = int($t1stMatchS2);
        $series2 = "use_2nd_discont";

        $use = "hitx2_useFIRST_tooSmall2ndHit";
        $temp_scoreBP = $scoreBP1;
        $temp_1stMatch = $t1stMatch1;
        $temp_strand = $strand1;
        $temp_1stMatchS = $t1stMatchS1;

        $scoreBP1 = $scoreBP1;
        $firstM1 = $t1stMatch1;
    }
}

```

```

$firstM2 = $tstMatch2;
$strand2 = $tstrand2;
$firstMS2 = int($tstMatchS2);
$series2 = "subset";

# (3.3): second hit begins first
} else {

    if ( (int($scoreBP2) >= $cf3 ) && (int($percent2) >= $cf2 ) ) {
        if ( (int($tstMatch2)+int($scoreBP2) >= (int($tstMatch1)+2) ) {
            $use = "hitx2_useSECONDnFIRST_continuous";
            # DoubleHPV-- independent of HPV16's 4 groups
            print DUB $Qnum,"#", $QL,"#hitHPV#","$use","\n";
            $temp_scoreBP = int($scoreBP1)+(int($tstMatch1)
                -int($tstMatch2));
            $temp_1stMatch = $tstMatch2;
            $temp_strand = $tstrand2;
            $scoreBP2 = $tstMatch2;
            $firstM2 = $tstMatch2;
            $strand2 = $tstrand2;

            $temp_1stMatchS = $tstMatchS2;
            $scoreBP1 = $tstMatchBP1;
            $firstM1 = $tstMatch1;
            $strand1 = $tstrand1;
            $firstMS1 = int($tstMatchS1);
            $series1 = "use_2nd";

            $firstMS2 = int($tstMatchS2);
            $series2 = "use_1st";

        } else {

            $use = "hitx2_useSECONDnFIRST_discontinuous";
            # DoubleHPV-- independent of HPV16's 4 groups
            print DUB $Qnum,"#", $QL,"#hitHPV#","$use","\n";
            $temp_scoreBP = $tstMatchBP2;
            $temp_1stMatch = $tstMatch2;
            $temp_strand = $tstrand2;
            $scoreBP2 = $tstMatch2;
            $firstM2 = $tstMatch2;
            $strand2 = $tstrand2;

            $temp_1stMatchS = $tstMatchS2;
            $scoreBP1 = $tstMatchBP1;
            $firstM1 = $tstMatch1;
            $strand1 = $tstrand1;
            $firstMS1 = int($tstMatchS1);
            $series1 = "use_2nd_discont";

            $firstMS2 = int($tstMatchS2);
            $series2 = "use_1st";

        }

        $temp_scoreBP = $tstMatchBP1;
        $temp_1stMatch = $tstMatch1;
        $temp_strand = $tstrand1;
        $use = "hitx2_useFIRST_tooSmall2ndHit";
        $temp_1stMatchS = $tstMatchS1;

        $scoreBP1 = $tstMatchBP1;
        $firstM1 = $tstMatch1;
        $strand1 = $tstrand1;
        $firstMS1 = int($tstMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $tstMatchBP2;
        $firstM2 = $tstMatch2;
        $strand2 = $tstrand2;
        $firstMS2 = int($tstMatchS2);
        $series2 = "smallerthan_",$cf3;

    }

    # for no-hit
} else {
    $use = "noHPV";
}

```

```

# (2.3): second hit begins first
} else {

    if ((int($percent1)>=int($percent2))&&(int($tstMatch1)-int($tstMatch2)<3)){
        $temp_scoreBP = $tstMatchBP1;
        $temp_1stMatch = $tstMatch1;
        $temp_strand = $tstrand1;
        $use = "hitx2_useFIRST_2ndHitasSubset";

        $temp_1stMatchS = $tstMatchS1;
        $scoreBP1 = $tstMatchBP1;
        $firstM1 = $tstMatch1;
        $strand1 = $tstrand1;
        $firstMS1 = int($tstMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $tstMatchBP2;
        $firstM2 = $tstMatch2;
        $strand2 = $tstrand2;
        $firstMS2 = int($tstMatchS2);
        $series2 = "subset";

        $temp_scoreBP = $tstMatchBP2;
        $temp_1stMatch = $tstMatch2;
        $temp_strand = $tstrand2;
        $use = "hitx2_useSECOND_1stHitasSubset";
        $temp_1stMatchS = $tstMatchS2;

        $scoreBP1 = $tstMatchBP1;
        $firstM1 = $tstMatch1;
        $strand1 = $tstrand1;
        $firstMS1 = int($tstMatchS1);
        $series1 = "subset";

        $scoreBP2 = $tstMatchBP2;
        $firstM2 = $tstMatch2;
        $strand2 = $tstrand2;
        $firstMS2 = int($tstMatchS2);
        $series2 = "use_Alonge";

    }

    # (3): if ending-of-first-hit is LATER than ending-of-second-hit
} else {

    # (3.1): first & second hits begin at the same pos
    if ( int($tstMatch1) == int($tstMatch2) ) {
        $temp_scoreBP = $tstMatchBP1;
        $temp_1stMatch = $tstMatch1;
        $temp_strand = $tstrand1;
        $use = "hitx2_useFIRST_2ndHitasSubset";
        $temp_1stMatchS = $tstMatchS1;

        $scoreBP1 = $tstMatchBP1;
        $firstM1 = $tstMatch1;
        $strand1 = $tstrand1;
        $firstMS1 = int($tstMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $tstMatchBP2;
        $firstM2 = $tstMatch2;
        $strand2 = $tstrand2;
        $firstMS2 = int($tstMatchS2);
        $series2 = "subset";

    }

    # (3.2): first hit begins first
    } elseif ( int($tstMatch1) < int($tstMatch2) ) {
        $temp_scoreBP = $tstMatchBP1;
        $temp_1stMatch = $tstMatch1;
        $temp_strand = $tstrand1;
        $use = "hitx2_useFIRST_2ndHitasSubset";
        $temp_1stMatchS = $tstMatchS1;

        $scoreBP1 = $tstMatchBP1;
        $firstM1 = $tstMatch1;
        $strand1 = $tstrand1;
        $firstMS1 = int($tstMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $tstMatchBP2;

```



```

File name:      set2_p5mac_d_28bpCutoff.pl
Source code:

#!/usr/bin/perl

## About the script #####
# Created by Sasithorn Chotewutmontri, Jan 2009.
#####

## MAIN PROGRAM BODY #####
#####
use strict;
use warnings;

my $sampleNo = @ARGV;
my $prefixPath = $sampleNo[2];
my $pyroNum = int($sampleNo[1]);

my $path = $prefixPath."Pyro".$pyroNum."_result/";
my $blastFilePath = $workingDir."blastHPV16/";

# For different cutoff conditions, change these variables:
#
my $newDir = $workingDir."28bpCutoff/";
my $prefix = "28bpCutoff_";

my $selectRep = $prefix."select_report";
my $grp1 = $prefix."select_group1";
my $grp2 = $prefix."select_group2";
my $grp3 = $prefix."select_group3";
my $grp4 = $prefix."select_group4";
my $selectReppub = $prefix."select_report_doubleHPV";

my $reportallpos = $prefix."allreport_fullpieceposition";

my $scf = 28;          # CUTOFF VALUE
my $scf2 = 94;         # CUTOFF PERCENT MATCH FOR 2nd HIT!!!!
my $scf3 = 15;         # CUTOFF SCOREBP FOR 2nd HIT!!!!

##### Start of PART 1 : SELECT into 4 groups #####
my $parsedFile = $blastFilePath."parsed_hpv_allreport";

my @allParseRep = ();

print "\n\nPROGRAM SD (SELECT 28bp CUTOFF) STARTS.....\n\n";
print "\n\n 4 'significantHPV16 hit' groups (using CF=28 nt)\n\n";
print "\n\n according to its 'non-HPV16' part\n\n";
print " G1 >= the non-HPV is >= 150 nt\n";
print " G2 >= the non-HPV is >= 100 AND < 150 nt\n";
print " G3 >= the non-HPV is >= 50 AND < 100 nt\n";
print " G4 >= the non-HPV is >= 0 AND < 50 nt\n";

# read parsed data for selection
@allParseRep = getfileData($parsedFile);

# make Dir for selected result of 'no cut-off'-program
System ("mkdir", $newDir);

unless (open (GRP1, ">$newDir$scfRep")) { print "\nCan not create ".$newDir.$scfRep."\n\n"; exit; }
unless (open (GRP2, ">$newDir$scf2Rep")) { print "\nCan not create ".$newDir.$scf2Rep."\n\n"; exit; }
unless (open (GRP3, ">$newDir$scf3Rep")) { print "\nCan not create ".$newDir.$scf3Rep."\n\n"; exit; }
unless (open (GRP4, ">$newDir$scf4Rep")) { print "\nCan not create ".$newDir.$scf4Rep."\n\n"; exit; }
unless (open (DUB, ">$newDir$selectReppub")) { print "\nCan not create ".$newDir.$selectReppub."\n\n"; exit; }
unless (open (RALL, ">$newDir$reportallpos")) { print "\nCan not create ".$newDir.$reportallpos."\n\n"; exit; }

for ( my $p = 0 ; $p < (scalar @allParseRep) ; $p++) {
    my $tmp_line = split (/./, $allParseRep[$p]);
    my $num = $tmp_line[0];
    my $SQL = $tmp_line[1];

    # query's seqNo
    # query's length

    my $tmp_index = $tmp_line[2];          # hit/nohit indice

    my $blockCount = scalar (@tmp_line);

    my ($scoreBP2, $firstM2, $strand2, $tmp_strand, $tmp_1stMatchS, $use);
    my ($scoreBP1, $firstM1, $strand1, $tmp_strand1, $tmp_1stMatchS1, $series1, $series2);

    # This first IF-loop defines the 'virtual matching area of HPV16'
    # The value of the 'first position' of HPV-matching area will be specified
    # as well as other values, e.g. 'virtual scoreBP' in case of two continuous hits

    # ONE-hit only
    if ((($tmp_index =~ "hitHPV") && ($blockCount == 4)) { # for ONE-HPV-HIT query
        my @tmpBox = split ( /%/, $tmp_line[3]);          # scoreBP
        $tmp_scoreBP = $tmpBox[0];
        $tmp_1stMatch = $tmpBox[2];
        $tmp_strand = $tmpBox[3];
        $tmp_1stMatchS = $tmpBox[4];
        $use = "hitx1_userIRST";

        $scoreBP1 = $tmpBox[0];
        $firstM1 = $tmpBox[2];
        $strand1 = $tmpBox[3];
        $firstM1S1 = $tmpBox[4];
        $series1 = "use_Alone";

        $scoreBP2 = "none";
        $firstM2 = "none";
        $strand2 = "none";
        $firstM2S2 = "none";
        $series2 = "none";

        # More than ONE hit, then two first hits will be taken into consideration
    } elsif (($tmp_index =~ "hitHPV") && ($blockCount > 4)) { # for MORE-HPV-HITS query
        my @tBox1 = split ( /%/, $tmp_line[3]);
        my $scoreBP1 = $tBox1[0];
        my $percent1 = $tBox1[1];
        my $1stMatch1 = $tBox1[2];
        my $strand1 = $tBox1[3];
        my $1stMatchS1 = $tBox1[4];

        my @tBox2 = split ( /%/, $tmp_line[4]);
        my $scoreBP2 = $tBox2[0];
        my $percent2 = $tBox2[1];
        my $1stMatch2 = $tBox2[2];
        my $strand2 = $tBox2[3];
        my $1stMatchS2 = $tBox2[4];

        # (1) : if ending-of-first-hit is EARLIER than ending-of-second-hit
        if ((int($1stMatch1)+int($scoreBP1)) < (int($1stMatch2)+int($scoreBP2))) {
            # (1.1): first & second hits begin at the same pos
            if ( int($1stMatch2) ) {
                $tmp_scoreBP = int($1stMatch2);
                $tmp_1stMatch = $1stMatch2;
                $tmp_strand = $strand2;
                $use = "hitx2_useSECOND_1stHitSubset";
                $tmp_1stMatchS = $1stMatchS2;

                $scoreBP1 = $1stMatchS1;
                $firstM1 = $1stMatch1;
                $strand1 = $strand1;
                $firstM1S1 = int($1stMatchS1);
                $series1 = "subset";

                $scoreBP2 = $1stMatchS2;
                $firstM2 = $1stMatch2;
                $strand2 = $strand2;
                $firstM2S2 = int($1stMatchS2);
                $series2 = "use_Alone";

                # (1.2): first hit begins first
            } elsif ( int($1stMatch1) < int($1stMatch2) ) {

```

```

$firstM1 = $tstMatch1;
$strand1 = $tstrand1;
$firstMS1 = int($tstMatchS1);
$series1 = "use_Alonge";

$scoreBP2 = $tscoreBP2;
$firstM2 = $tstMatch2;
$strand2 = $tstrand2;
$firstMS2 = int($tstMatchS2);
$series2 = "smallerThan_".scf3;
}
}

# (1.3): second hit begins first
} else {
    $temp_scoreBP = $tscoreBP2;
    $temp_1stMatch = $tstMatch2;
    $temp_strand = $tstrand2;
    $use = "hitx2_useSECOND_1stHitAsSubset";
    $temp_1stMatchS = $tstMatchS2;

    $scoreBP1 = $tscoreBP1;
    $firstM1 = $tstMatch1;
    $strand1 = $tstrand1;
    $firstMS1 = int($tstMatchS1);
    $series1 = "subset";

    $scoreBP2 = $tscoreBP2;
    $firstM2 = $tstMatch2;
    $strand2 = $tstrand2;
    $firstMS2 = int($tstMatchS2);
    $series2 = "use_Alonge";
}

# (2) : if ending-of-first-hit is THE SAME as ending-of-second-hit
} elseif ((int($tstMatch1)+int($tscoreBP1))==(int($tstMatch2)+int($tscoreBP2))) {
    # (2.1): first & second hits begin at the same pos
    if ( int($tstMatch1) == int($tstMatch2) ) {
        $temp_scoreBP = $tscoreBP1;
        $temp_1stMatch = $tstMatch1;
        $temp_strand = $tstrand1;
        $use = "hitx2_useFIRST_2ndHitAsSubset";
        $temp_1stMatchS = $tstMatchS1;

        $scoreBP1 = $tscoreBP1;
        $firstM1 = $tstMatch1;
        $strand1 = $tstrand1;
        $firstMS1 = int($tstMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $tscoreBP2;
        $firstM2 = $tstMatch2;
        $strand2 = $tstrand2;
        $firstMS2 = int($tstMatchS2);
        $series2 = "subset";
    }

    # (2.2): first hit begins first
    } elseif ( int($tstMatch1) < int($tstMatch2) ) {
        $temp_scoreBP = $tscoreBP1;
        $temp_1stMatch = $tstMatch1;
        $temp_strand = $tstrand1;
        $use = "hitx2_useFIRST_2ndHitAsSubset";
        $temp_1stMatchS = $tstMatchS1;

        $scoreBP1 = $tscoreBP1;
        $firstM1 = $tstMatch1;
        $strand1 = $tstrand1;
        $firstMS1 = int($tstMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $tscoreBP2;
        $firstM2 = $tstMatch2;
        $strand2 = $tstrand2;
        $firstMS2 = int($tstMatchS2);
        $series2 = "subset";
    }

    # (2.3): second hit begins first
    } else {

```

```

# If last match position of 1stHit is >= 1stMatch of 2.Hit (+2)
# (+2) is the allowance for the gap or if there is no gap then 2 nt distance
# between the two hits
if ( (int($tstMatch1)+int($tscoreBP1)) >= (int($tstMatch2)+2) ) {

    # Control if the 2.Hit is long enough to be significantly counted
    if (((int($tscoreBP2) >= scf3 ) && (int($percent2) >= scf2 )) {
        $use = "hitx2_useFIRSTnSECOND_continuous";
        # DoubleHPV-- independent of HPV16's 4 groups

        print DUB $Qnum."#".SQL."#hitHPV#". $use."\\n";
        $temp_scoreBP = int($tscoreBP2)+(int($tstMatch2)
            -int($tstMatch1));

        $temp_1stMatch = $tstMatch1;
        $temp_strand = $tstrand1;
        $scoreBP2 = $tscoreBP2;
        $firstM2 = $tstMatch2;
        $strand2 = $tstrand2;
        $temp_1stMatchS = $tstMatchS1;

        $scoreBP1 = $tscoreBP1;
        $firstM1 = $tstMatch1;
        $strand1 = $tstrand1;
        $firstMS1 = int($tstMatchS1);
        $series1 = "use_1st";

        $firstMS2 = int($tstMatchS2);
        $series2 = "use_2nd";

        $use = "hitx2_useFIRST_tooSmall2ndHit";
        $temp_scoreBP = $tscoreBP1;
        $temp_1stMatch = $tstMatch1;
        $temp_strand = $tstrand1;
        $temp_1stMatchS = $tstMatchS1;

        $scoreBP1 = $tscoreBP1;
        $firstM1 = $tstMatch1;
        $strand1 = $tstrand1;
        $firstMS1 = int($tstMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $tscoreBP2;
        $firstM2 = $tstMatch2;
        $strand2 = $tstrand2;
        $firstMS2 = int($tstMatchS2);
        $series2 = "smallerThan_".scf3;
    }

    } else {
        if ( (int($tscoreBP2) >= scf3 ) && (int($percent2) >= scf2 )) {
            $use = "hitx2_useFIRSTnSECOND_discontinuous";
            # DoubleHPV-- independent of HPV16's 4 groups

            print DUB $Qnum."#".SQL."#hitHPV#". $use."\\n";
            $temp_scoreBP = $tscoreBP1; # DIFFERENT FROM "continuous"
            $temp_1stMatch = $tstMatch1;
            $temp_strand = $tstrand1;
            $scoreBP2 = $tscoreBP2;
            $firstM2 = $tstMatch2;
            $strand2 = $tstrand2;
            $temp_1stMatchS = $tstMatchS1;

            $scoreBP1 = $tscoreBP1;
            $firstM1 = $tstMatch1;
            $strand1 = $tstrand1;
            $firstMS1 = int($tstMatchS1);
            $series1 = "use_1st";

            $firstMS2 = int($tstMatchS2);
            $series2 = "use_2nd_discont";
        }

        $use = "hitx2_useFIRST_tooSmall2ndHit";
        $temp_scoreBP = $tscoreBP1;
        $temp_1stMatch = $tstMatch1;
        $temp_strand = $tstrand1;
        $temp_1stMatchS = $tstMatchS1;

        $scoreBP1 = $tscoreBP1;

```

```

$firstMS2 = int($t1stMatchS2);
$series2 = "subset";

# (3.3): second hit begins first
} else {
    if ( (int($scoreBP2) >= $cf3 ) && (int($percent2) >= $cf2 ) ) {
        if ( (int($t1stMatch2)+int($scoreBP2)) >= (int($t1stMatch1)+2) ) {
            $use = "hitx2_useSECONDnFIRST_continuous";
            # DoubleHPV-- independent of HPV16's 4 groups
            print DUB $Qnum,"#", $QL,"#hitHPV#",$use,"\\n";
            $temp_scoreBP = int($scoreBP1)+(int($t1stMatch1
                -int($t1stMatch2)));
            $temp_1stMatch = $t1stMatch2;
            $temp_strand = $t1strand2;
            $scoreBP2 = $t1scoreBP2;
            $firstM2 = $t1stMatch2;
            $strand2 = $t1strand2;

            $temp_1stMatchS = $t1stMatchS2;
            $scoreBP1 = $t1scoreBP1;
            $firstM1 = $t1stMatch1;
            $strand1 = $t1strand1;
            $firstMS1 = int($t1stMatchS1);
            $series1 = "use_2nd";

            $firstMS2 = int($t1stMatchS2);
            $series2 = "use_1st";

            $use = "hitx2_useSECONDnFIRST_discontinuous";
            # DoubleHPV-- independent of HPV16's 4 groups
            print DUB $Qnum,"#", $QL,"#hitHPV#",$use,"\\n";
            $temp_scoreBP = $t1scoreBP2;
            $temp_1stMatch = $t1stMatch2;
            $temp_strand = $t1strand2;
            $scoreBP2 = $t1scoreBP2;
            $firstM2 = $t1stMatch2;
            $strand2 = $t1strand2;

            $temp_1stMatchS = $t1stMatchS2;
            $scoreBP1 = $t1scoreBP1;
            $firstM1 = $t1stMatch1;
            $strand1 = $t1strand1;
            $firstMS1 = int($t1stMatchS1);
            $series1 = "use_2nd_discont";

            $firstMS2 = int($t1stMatchS2);
            $series2 = "use_1st";
        }
    }
    $temp_scoreBP = $t1scoreBP1;
    $temp_1stMatch = $t1stMatch1;
    $temp_strand = $t1strand1;
    $use = "hitx2_useFIRST_tooSmall2ndHit";
    $temp_1stMatchS = $t1stMatchS1;
    $scoreBP2 = $t1scoreBP1;
    $firstM1 = $t1stMatch1;
    $strand1 = $t1strand1;
    $firstMS1 = int($t1stMatchS1);
    $series1 = "use_Alonge";

    $scoreBP2 = $t1scoreBP2;
    $firstM2 = $t1stMatch2;
    $strand2 = $t1strand2;
    $firstMS2 = int($t1stMatchS2);
    $series2 = "smallerthan_".$cf3;
}

}

} else { # for no-hit
    $use = "noHPV";
}
}

```

```

if ((int($percent1)>=int($percent2))&&((int($t1stMatch1)-int($t1stMatch2))<3) ) {
    $temp_scoreBP = $t1scoreBP1;
    $temp_1stMatch = $t1stMatch1;
    $temp_strand = $t1strand1;
    $use = "hitx2_useFIRST_2ndHitasSubset";

    $temp_1stMatchS = $t1stMatchS1;
    $scoreBP1 = $t1scoreBP1;
    $firstM1 = $t1stMatch1;
    $strand1 = $t1strand1;
    $firstMS1 = int($t1stMatchS1);
    $series1 = "use_Alonge";

    $scoreBP2 = $t1scoreBP2;
    $firstM2 = $t1stMatch2;
    $strand2 = $t1strand2;
    $firstMS2 = int($t1stMatchS2);
    $series2 = "subset";

    $temp_scoreBP = $t1scoreBP2;
    $temp_1stMatch = $t1stMatch2;
    $temp_strand = $t1strand2;
    $use = "hitx2_useSECOND_1stHitasSubset";
    $temp_1stMatchS = $t1stMatchS2;

    $scoreBP1 = $t1scoreBP1;
    $firstM1 = $t1stMatch1;
    $strand1 = $t1strand1;
    $firstMS1 = int($t1stMatchS1);
    $series1 = "subset";

    $scoreBP2 = $t1scoreBP2;
    $firstM2 = $t1stMatch2;
    $strand2 = $t1strand2;
    $firstMS2 = int($t1stMatchS2);
    $series2 = "use_Alonge";
}

# (3) : if ending-of-first-hit is LATER than ending-of-second-hit
} else {
    # (3.1): first & second hits begin at the same pos
    if ( int($t1stMatch1) == int($t1stMatch2) ) {
        $temp_scoreBP = $t1scoreBP1;
        $temp_1stMatch = $t1stMatch1;
        $temp_strand = $t1strand1;
        $use = "hitx2_useFIRST_2ndHitasSubset";
        $temp_1stMatchS = $t1stMatchS1;

        $scoreBP1 = $t1scoreBP1;
        $firstM1 = $t1stMatch1;
        $strand1 = $t1strand1;
        $firstMS1 = int($t1stMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $t1scoreBP2;
        $firstM2 = $t1stMatch2;
        $strand2 = $t1strand2;
        $firstMS2 = int($t1stMatchS2);
        $series2 = "subset";

        # (3.2): first hit begins first
    } elseif ( int($t1stMatch1) < int($t1stMatch2) ) {
        $temp_scoreBP = $t1scoreBP1;
        $temp_1stMatch = $t1stMatch1;
        $temp_strand = $t1strand1;
        $use = "hitx2_useFIRST_2ndHitasSubset";
        $temp_1stMatchS = $t1stMatchS1;

        $scoreBP1 = $t1scoreBP1;
        $firstM1 = $t1stMatch1;
        $strand1 = $t1strand1;
        $firstMS1 = int($t1stMatchS1);
        $series1 = "use_Alonge";

        $scoreBP2 = $t1scoreBP2;
        $firstM2 = $t1stMatch2;
        $strand2 = $t1strand2;
    }
}

```


[illegible]

```
}
if (( $status == m/noHPV(.*$/) ) { ($status =~ m/no_group(.*/)) { $nohpv++; }
} else if (($status =~ m/G1(.*)$/) { ($G1 = $noG1b+); }
else if (($status =~ m/G2(.*)$/) { ($G2 = $noG2b+); }
else if (($status =~ m/G3(.*)$/) { ($G3 = $noG3b+); }
else if (($status =~ m/G4(.*)$/) { ($G4 = $noG4b+); }
else if (($status =~ m/out_ranged_G1.*)/$) { ($outG1++); }
else if (($status =~ m/out_ranged_G2.*)/$) { ($outG2++); }
else if (($status =~ m/out_ranged_G3.*)/$) { ($outG3++); }
else if (($status =~ m/out_ranged_G4.*)/$) { ($outG4++); }
}
my $noTotalSeq = (scalar @seqNameFile);
print OUT "Total number of sequences having this barcode is : ".$noTotalSeq."\n\n";
print OUT "CUT-OFF value used for minimum matching length is : ".$cutoff." bp.\n\n";
unless ($error != 0) {
    my $noUnclass = $noTotalSeq - $noG1 - $noG2 - $noG3 - $noG4;
    print OUT "Number of sequences in group 1 : ".$noG1."\n\n";
    print OUT "Number of sequences in group 2 : ".$noG2."\n\n";
    print OUT "Number of sequences in group 3 : ".$noG3."\n\n";
    print OUT "Number of sequences in group 4 : ".$noG4."\n\n";
    print OUT "Number of sequences with UNCLASSIFIED-status : ".$noUnclass."\n\n";
    print OUT " --> OUT-Ranged G1 : ".$outG1."\n\n";
    print OUT " --&#x2D;> OUT-Ranged G2 : ".$outG2."</code>.\\n\\n";
<code>    print OUT " --> OUT-Ranged G3 : ".$outG3.</code>"\\n\\n";
    print OUT " --> OUT-Ranged G4 : ".$outG4."</code>"\\n\\n";
    print OUT " --> no-hit-to-HPV : ".$nohpv."</code>"\\n\\n";
}
print OUT "\\nNote that UNCLASSIFIED can mean either \\n";
print OUT "a) OUT-RANGED from the four groups, i.e. Matching is shorter than CUTOFF";
print OUT "nor b) the sequence has NO-HIT-to-HPV\\n";
close OUT;
}

##### Sub-routines END HERE #####
##### SUBROUTINES #####
sub getRevConFastaG1to4 {
    my ($workDir, $fileName, $prefix) = @_;
    my $outdir = "$workDir.$prefix.out";
    my $fastaRC = "$workDir.$prefix.fastaRC";
    my %fas = getFileData($fastaRC);
    for (my $g = 1; $g <= 4; $g++) {
        my $groupList = $dir.$prefix."_select_group".$g.".names";
        my @list = getFileData($groupList);
        my $out = $dir.$prefix."_select_group".$g.".outsuffix";
        unless (open (FA, ">$out")) { print "\ncan't open file \"$out\""; exit; }
        for (my $sa = 0; $sa < scalar @list; $sa++) {
            my @box = split (/#/,$list[$sa]);
            my $seqnum = $box[0];
            if ($seqnum =~ /([0-9]*)$/) {
                my $sum = int($seqnum);
                my $infaName = $fastID{$sum};
                my $infaPath = $path.$infaNo{0}."/";
                my $infaSeq = $fastA{$infaNo{0}}+$1;
                print FA "$infaName\t"$infaSeq;
            }
        }
        close FA;
    }
}
##### End of Sub-routines #####
```



```

    } else if ($grp =~ /\out_ranged(G[0-9]S/) {
        $seOut++; print RE "num\tsspl[4]\tscate\tsgrp\n";
    } else if ($grp =~ /\no_groups/) { $seOut++; print RE "num\tsspl[4]\tscate\tsgrp\n"; }
    } else if ($grp =~ /\noHPS/) { $seIn++; print RE "num\tsspl[4]\tscate\tsgrp\n"; }
    } else { print "\nError\n\n"; exit; }
} elsif (($cate =~ /\OUT(.*)S/) {
    if ($grp =~ /\G1S/) { $oe5q1++; print RE "num\tinsig_sspl[4]\tscate\tsgrp\n"; }
    if ($grp =~ /\G2S/) { $oe5q2++; print RE "num\tinsig_sspl[4]\tscate\tsgrp\n"; }
    if ($grp =~ /\G3S/) { $oe5q3++; print RE "num\tinsig_sspl[4]\tscate\tsgrp\n"; }
    if ($grp =~ /\G4S/) { $oe5q4++; print RE "num\tinsig_sspl[4]\tscate\tsgrp\n"; }
    } else if ($grp =~ /\out_ranged(G[0-9]S/) {
        $seOut++; print RE "num\tinsig_sspl[4]\tscate\tsgrp\n"; }
    } else if ($grp =~ /\no_groups/) {
        $seOut++; print RE "num\tinsig_sspl[4]\tscate\tsgrp\n"; }
    } else if ($grp =~ /\noHPS/) { $oe5nh++; print RE "num\tinsig_sspl[4]\tscate\tsgrp\n"; }
    } else { print "\nError\n\n"; exit; }
} else {
    print "\n\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;
}
}
} elsif ($prtm =~ /\$spl[5]S/) {
    if (($cate =~ /\completes/) || ($cate =~ /\INcom(.*)S/) ) {
        if ($grp =~ /\G1S/) {
            $se6q1++; print RE "num\tsspl[5]\tscate\tsgrp\n"; print LE6 " $num\n\n";
        }
        $se6q2++; print RE "num\tsspl[5]\tscate\tsgrp\n"; print LE6 " $num\n\n";
        $se6q3++; print RE "num\tsspl[5]\tscate\tsgrp\n"; print LE6 " $num\n\n";
        $se6q4++; print RE "num\tsspl[5]\tscate\tsgrp\n"; print LE6 " $num\n\n";
        $se6q5++; print RE "num\tsspl[5]\tscate\tsgrp\n"; print LE6 " $num\n\n";
        } else if ($grp =~ /\out_ranged(G[0-9]S/) {
            $seOut++; print RE "num\tsspl[5]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\no_groups/) { $seOut++; print RE "num\tsspl[5]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\noHPS/) { $oe6nh++; print RE "num\tsspl[5]\tscate\tsgrp\n"; }
        } else { print "\nError\n\n"; exit; }
    } } elsif (($cate =~ /\OUT(.*)S/) {
        if ($grp =~ /\G1S/) { $oe6g1++; print RE "num\tinsig_sspl[5]\tscate\tsgrp\n"; }
        if ($grp =~ /\G2S/) { $oe6g2++; print RE "num\tinsig_sspl[5]\tscate\tsgrp\n"; }
        if ($grp =~ /\G3S/) { $oe6g3++; print RE "num\tinsig_sspl[5]\tscate\tsgrp\n"; }
        if ($grp =~ /\G4S/) { $oe6g4++; print RE "num\tinsig_sspl[5]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\out_ranged(G[0-9]S/) {
            $seOut++; print RE "num\tinsig_sspl[5]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\no_groups/) {
            $seOut++; print RE "num\tinsig_sspl[5]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\noHPS/) { $oe6nh++; print RE "num\tinsig_sspl[5]\tscate\tsgrp\n"; }
        } else { print "\nError\n\n"; exit; }
    } } else {
        print "\n\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;
    }
}
} elsif ($prtm =~ /\$spl[6]S/) {
    if (($cate =~ /\completes/) || ($cate =~ /\INcom(.*)S/) ) {
        if ($grp =~ /\G1S/) {
            $se7q1++; print RE "num\tsspl[6]\tscate\tsgrp\n"; print LE7 " $num\n\n";
        }
        $se7q2++; print RE "num\tsspl[6]\tscate\tsgrp\n"; print LE7 " $num\n\n";
        $se7q3++; print RE "num\tsspl[6]\tscate\tsgrp\n"; print LE7 " $num\n\n";
        $se7q4++; print RE "num\tsspl[6]\tscate\tsgrp\n"; print LE7 " $num\n\n";
        $se7q5++; print RE "num\tsspl[6]\tscate\tsgrp\n"; print LE7 " $num\n\n";
        } else if ($grp =~ /\out_ranged(G[0-9]S/) {
            $seOut++; print RE "num\tsspl[6]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\no_groups/) { $seOut++; print RE "num\tsspl[6]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\noHPS/) { $se7nh++; print RE "num\tsspl[6]\tscate\tsgrp\n"; }
        } else { print "\nError\n\n"; exit; }
    } } elsif (($cate =~ /\OUT(.*)S/) {
        if ($grp =~ /\G1S/) { $oe7g1++; print RE "num\tinsig_sspl[6]\tscate\tsgrp\n"; }
        if ($grp =~ /\G2S/) { $oe7g2++; print RE "num\tinsig_sspl[6]\tscate\tsgrp\n"; }
        if ($grp =~ /\G3S/) { $oe7g3++; print RE "num\tinsig_sspl[6]\tscate\tsgrp\n"; }
        if ($grp =~ /\G4S/) { $oe7g4++; print RE "num\tinsig_sspl[6]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\out_ranged(G[0-9]S/) {
            $seOut++; print RE "num\tinsig_sspl[6]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\no_groups/) {
            $seOut++; print RE "num\tinsig_sspl[6]\tscate\tsgrp\n"; }
        } else if ($grp =~ /\noHPS/) { $oe7nh++; print RE "num\tinsig_sspl[6]\tscate\tsgrp\n"; }
        } else { print "\nError\n\n"; exit; }
    } } else {
        print "\n\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;
    }
}
} elsif ($prtm =~ /\$spl[7]S/) {

```

[illegible]

[illegible]

[illegible]

```

    elseif ($grp ~ /^no_groups/) { $sel3out++; print RE "$num\t$sspl[12]\t$cate\t$sgprv\n";  
    elseif ($grp =~ /noHPVs/) { $sel3nh++; print RE "$num\t$sspl[12]\t$cate\t$sgprv\n";  
    else {print "\nError\n\n"; exit;}  
}  
} elseif (($cate =~ /^OUT(.*)$/) ) {  
    if ($grp =~ ^/G1S/) { $sel3opl++; print RE " $num\t$nsig_ $spl[12]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/G2S/) { $sel3o2++; print RE " $num\t$nsig_ $spl[12]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/G3S/) { $sel3o3++; print RE " $num\t$nsig_ $spl[12]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/G4S/) { $sel3o4++; print RE " $num\t$nsig_ $spl[12]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/out_ranged_G(0-9)s/) {  
        $sel3out+;; print RE "$num\t$nsig_ $spl[12]\t$cate\t$sgprv\n";  
    }  
    elseif ($grp =~ /^no_groups/) {  
        $sel3out+;; print RE " $num\t$nsig_ $spl[12]\t$cate\t$sgprv\n";  
    }  
    elseif ($grp =~ /noHPVs/) { $sel3nh++; print RE "$num\t$nsig_ $spl[12]\t$cate\t$sgprv\n";  
    else {print "\nError\n\n"; exit;}  
}  
} else {  
    print "\n\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;  
}  
}  
} elseif ($prim =~ /^$spl[13]/) {  
    if (($cate =~ /^completes/) || ($cate =~ /^TnCom(.*)$/)) {  
        if ($grp =~ ^/G1S/) {  
            $sel4q1+; print RE " $num\t$sspl[13]\t$cate\t$gprv\n"; print LL6 "$num\n\n";  
            $sel4q2+; print RE " $num\t$sspl[13]\t$cate\t$gprv\n"; print LL6 "$num\n\n";  
            $sel4q3+; print RE " $num\t$sspl[13]\t$cate\t$gprv\n"; print LL6 "$num\n\n";  
            $sel4q4+; print RE " $num\t$sspl[13]\t$cate\t$gprv\n"; print LL6 "$num\n\n";  
            $sel4q5+; print RE " $num\t$sspl[13]\t$cate\t$gprv\n"; print LL6 "$num\n\n";  
            $sel4out+; print RE " $num\t$sspl[13]\t$cate\t$sgprv\n";  
        }  
    }  
    elseif ($grp =~ /^no_groups/) { $sel4out++; print RE " $num\t$sspl[13]\t$cate\t$sgprv\n";  
    elseif ($grp =~ /noHPVs/) { $sel4nh++; print RE " $num\t$sspl[13]\t$cate\t$sgprv\n";  
    else {print "\nError\n\n"; exit;}  
}  
} elseif (($cate =~ /^OUT(.*)$/) ) {  
    if ($grp =~ ^/G1S/) { $sel4o1++; print RE " $num\t$nsig_ $spl[13]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/G2S/) { $sel4o2++; print RE " $num\t$nsig_ $spl[13]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/G3S/) { $sel4o3++; print RE " $num\t$nsig_ $spl[13]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/G4S/) { $sel4o4++; print RE " $num\t$nsig_ $spl[13]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/out_ranged_G(0-9)s/) {  
        $sel4out+; print RE " $num\t$nsig_ $spl[13]\t$cate\t$sgprv\n";  
    }  
    elseif ($grp =~ /^no_groups/) {  
        $sel4out+; print RE " $num\t$nsig_ $spl[13]\t$cate\t$sgprv\n";  
    }  
    elseif ($grp =~ /noHPVs/) { $sel4nh++; print RE " $num\t$nsig_ $spl[13]\t$cate\t$sgprv\n";  
    else {print "\nError\n\n"; exit;}  
}  
} else {  
    print "\n\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;  
}  
}  
} elseif ($prim =~ /^$spl[14]/) {  
    if (($cate =~ /^completes/) || ($cate =~ /^TnCom(.*)$/)) {  
        if ($grp =~ ^/G1S/) {  
            $sel5q1+; print RE " $num\t$sspl[14]\t$cate\t$gprv\n"; print LL7 "$num\n\n";  
            $sel5q2+; print RE " $num\t$sspl[14]\t$cate\t$gprv\n"; print LL7 "$num\n\n";  
            $sel5q3+; print RE " $num\t$sspl[14]\t$cate\t$gprv\n"; print LL7 "$num\n\n";  
            $sel5q4+; print RE " $num\t$sspl[14]\t$cate\t$gprv\n"; print LL7 "$num\n\n";  
            $sel5out+; print RE " $num\t$sspl[14]\t$cate\t$sgprv\n"; print LL7 "$num\n\n";  
            $sel5o1+; print RE " $num\t$sspl[14]\t$cate\t$sgprv\n";  
        }  
    }  
    elseif ($grp =~ /^no_groups/) { $sel5out++; print RE " $num\t$sspl[14]\t$cate\t$sgprv\n";  
    elseif ($grp =~ /noHPVs/) { $sel5nh++; print RE " $num\t$sspl[14]\t$cate\t$sgprv\n";  
    else {print "\nError\n\n"; exit;}  
}  
} elseif (($cate =~ /^OUT(.*)$/) ) {  
    if ($grp =~ ^/G1S/) { $sel5o1++; print RE " $num\t$nsig_ $spl[14]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/G2S/) { $sel5o2++; print RE " $num\t$nsig_ $spl[14]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/G3S/) { $sel5o3++; print RE " $num\t$nsig_ $spl[14]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/G4S/) { $sel5o4++; print RE " $num\t$nsig_ $spl[14]\t$cate\t$sgprv\n";  
    elseif ($grp =~ ^/out_ranged_G(0-9)s/) {  
        $sel5out+; print RE " $num\t$nsig_ $spl[14]\t$cate\t$sgprv\n";  
    }  
    elseif ($grp =~ /^no_groups/) {  
        $sel5out+; print RE " $num\t$nsig_ $spl[14]\t$cate\t$sgprv\n";  
    }  
    elseif ($grp =~ /noHPVs/) { $sel5nh++; print RE " $num\t$nsig_ $spl[14]\t$cate\t$sgprv\n";  
    else {print "\nError\n\n"; exit;}  
}  
} else {  
    print "\n\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;  
}  
}  
} elseif ($prim =~ /^$spl[15]/) {  
    if (($cate =~ /^completes/) || ($cate =~ /^TnCom(.*)$/)) {  
        if ($grp =~ ^/G1S/) {  


```

```

                se$23out+; print RE "$num\tsspl [20]\tscate\tsgrp\n";}
            elseif (se$rp ~ /noHPVs/) {se$23nh+; print RE "$num\tsspl [20]\tscate\tsgrp\n";}
            else {print "\nError!\n\n"; exit;}
        }
    } elseif ((scate == "/OUT(.*?)$/")) {
        if ($rp == ~/G1S/) {$se$23oi+; print RE "num\tinsig_ $spl [20]\tscate\tsgrp\n";}
        elseif ($rp == ~/G2S/) {$se$23oi+; print RE "num\tinsig_ $spl [20]\tscate\tsgrp\n";}
        elseif ($rp == ~/G3S/) {$se$23oi+; print RE "num\tinsig_ $spl [20]\tscate\tsgrp\n";}
        elseif ($rp == ~/G4S/) {$se$23oi+; print RE "num\tinsig_ $spl [20]\tscate\tsgrp\n";}
        elseif ($rp == /out_ranged_G([0-9])$/) {
            se$23out+; print RE "$num\tinsig_ $spl [20]\tscate\tsgrp\n";}
        elseif ($rp == /no_groups/) {
            se$23out+; print RE "num\tinsig_ $spl [20]\tscate\tsgrp\n";}
        elseif ($rp == /noHPVs/) {$se$23nh+; print RE "$num\tinsig_ $spl [20]\tscate\tsgrp\n";}
        else {print "\nError!\n\n"; exit;}
    }
} else {
    print "\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;
}

# elif ($prm =~ /\$spl[22]$/) {
    If (($scate =~ /^completes/) || ($cate =~ /^INcom(.*)$/)) {
        if ($rp == ~/G1S/) {
            se$24oi+; print RE "$num\tsspl [21]\tscate\tsgrp\n"; print LE24 "$num\n";}
            elseif ($rp == ~/G2S/) {
            se$24oi+; print RE "$num\tsspl [21]\tscate\tsgrp\n"; print LE24 "$num\n";}
            elseif ($rp == ~/G3S/) {
            se$24oi+; print RE "$num\tsspl [21]\tscate\tsgrp\n"; print LE24 "$num\n";}
            elseif ($rp == ~/G4S/) {
            se$24oi+; print RE "$num\tsspl [21]\tscate\tsgrp\n"; print LE24 "$num\n";}
            elseif ($rp == /out_ranged_G([0-9])$/) {
            se$24oi+; print RE "num\tsspl [21]\tscate\tsgrp\n";}
            elseif ($rp == /no_groups/) {
            se$24oi+; print RE "num\tsspl [21]\tscate\tsgrp\n";}
            elseif ($rp == ~/noHPVs/) {$se$24out+; print RE "num\tsspl [21]\tscate\tsgrp\n";}
            else {print "\nError!\n\n"; exit;}
        }
    } elseif ((scate =~ "/OUT(.*?)$/")) {
        if ($rp == ~/G1S/) {$se$24oi+; print RE "num\tinsig_ $spl [21]\tscate\tsgrp\n";}
        elseif ($rp == ~/G2S/) {$se$24oi+; print RE "num\tinsig_ $spl [21]\tscate\tsgrp\n";}
        elseif ($rp == ~/G3S/) {$se$24oi+; print RE "num\tinsig_ $spl [21]\tscate\tsgrp\n";}
        elseif ($rp == ~/G4S/) {$se$24oi+; print RE "num\tinsig_ $spl [21]\tscate\tsgrp\n";}
        elseif ($rp == /out_ranged_G([0-9])$/) {
            se$24out+; print RE "num\tinsig_ $spl [21]\tscate\tsgrp\n";}
        elseif ($rp == /no_groups/) {
            se$24out+; print RE "num\tinsig_ $spl [21]\tscate\tsgrp\n";}
        elseif ($rp == /noHPVs/) {$se$24nh+; print RE "num\tinsig_ $spl [21]\tscate\tsgrp\n";}
        else {print "\nError!\n\n"; exit;}
    }
} else {
    print "\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;
}

# elif ($prm =~ /\$spl[22]$/) {
    If (($scate =~ /^completes/) || ($cate =~ /^INcom(.*)$/)) {
        if ($rp == ~/G1S/) {
            se$25oi+; print RE "$num\tsspl [22]\tscate\tsgrp\n"; print LE25 "$num\n";}
            elseif ($rp == ~/G2S/) {
            se$25oi+; print RE "$num\tsspl [22]\tscate\tsgrp\n"; print LE25 "$num\n";}
            elseif ($rp == ~/G3S/) {
            se$25oi+; print RE "$num\tsspl [22]\tscate\tsgrp\n"; print LE25 "$num\n";}
            elseif ($rp == ~/G4S/) {
            se$25oi+; print RE "$num\tsspl [22]\tscate\tsgrp\n"; print LE25 "$num\n";}
            elseif ($rp == /out_ranged_G([0-9])$/) {
            se$25oi+; print RE "num\tsspl [22]\tscate\tsgrp\n";}
            elseif ($rp == /no_groups/) {
            se$25oi+; print RE "num\tsspl [22]\tscate\tsgrp\n";}
            elseif ($rp == ~/noHPVs/) {$se$25out+; print RE "num\tsspl [22]\tscate\tsgrp\n";}
            else {print "\nError!\n\n"; exit;}
        }
    } elseif ((scate =~ "/OUT(.*?)$/")) {
        if ($rp == ~/G1S/) {$se$25oi+; print RE "num\tinsig_ $spl [22]\tscate\tsgrp\n";}
        elseif ($rp == ~/G2S/) {$se$25oi+; print RE "num\tinsig_ $spl [22]\tscate\tsgrp\n";}
        elseif ($rp == ~/G3S/) {$se$25oi+; print RE "num\tinsig_ $spl [22]\tscate\tsgrp\n";}
        elseif ($rp == ~/G4S/) {$se$25oi+; print RE "num\tinsig_ $spl [22]\tscate\tsgrp\n";}
        elseif ($rp == /out_ranged_G([0-9])$/) {
            se$25out+; print RE "num\tinsig_ $spl [22]\tscate\tsgrp\n";}
        elseif ($rp == /no_groups/) {
            se$25out+; print RE "num\tinsig_ $spl [22]\tscate\tsgrp\n";}
        elseif ($rp == /noHPVs/) {$se$25nh+; print RE "num\tinsig_ $spl [22]\tscate\tsgrp\n";}
        else {print "\nError!\n\n"; exit;}
    }
} else {
    print "\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;
}

# elif ($prm =~ /\$spl[23]$/) {
    If (($scate =~ /^completes/) || ($cate =~ /^INcom(.*)$/)) {
        if ($rp == ~/G1S/) {

```

```

seoe19out++; print RE "$num\tINSig_$$_$pl[17]\tscate\tsgrp\n"; }
elseif ($grp =~ /\nohpvs/) { $oe19inh++; print RE "$num\tINSig_$$_$pl[17]\tscate\tsgrp\n"; }
else { print "\nError\n\n"; exit; }
print "\n\nThere is error in PrimerMatchingCategory...$num\n\n";
exit;
} else {
    $prim =~ /\$spl[18]$/;
    if (($scate =~ /\completes/){($scate =~ /\INcom.(.*)$/)}{
        if ($grp =~ /\G1$/){
            se20g1+--; print RE "$num\t$$$pl[18]\tscate\tsgrp\n"; print LE20 "$num\n"; }
        elseif ($grp =~ /\G2$/){
            se20g2+--; print RE "$num\t$$$pl[18]\tscate\tsgrp\n"; print LE20 "$num\n"; }
        elseif ($grp =~ /\G3$/){
            se20g3+--; print RE "$num\t$$$pl[18]\tscate\tsgrp\n"; print LE20 "$num\n"; }
        elseif ($grp =~ /\G4$/){
            se20g4+--; print RE "$num\t$$$pl[18]\tscate\tsgrp\n"; print LE20 "$num\n"; }
        elseif ($grp =~ /\out_ranged.G(0-9)$/){
            se20out++; print RE "$num\t$$$pl[18]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\no.groups/) {
            se20out++; print RE "$num\t$$$pl[18]\tscate\tsgrp\n"; }
        else { print "\nError\n\n"; exit; }
    } elseif (($scate =~ /\OUT(.*)$/){
        if ($grp =~ /\G1$/){($oe20g1++); print RE "$num\tINSig_$$_$pl[18]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\G2$/){($oe20g2++); print RE "$num\tINSig_$$_$pl[18]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\G3$/){($oe20g3++); print RE "$num\tINSig_$$_$pl[18]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\G4$/){($oe20g4++); print RE "$num\tINSig_$$_$pl[18]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\out_ranged.G(0-9)$/){
            se20out++; print RE "$num\tINSig_$$_$pl[18]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\no.groups/) {
            se20out++; print RE "$num\tINSig_$$_$pl[18]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\noHPVs/) {($oe20nh++); print RE "$num\tINSig_$$_$pl[18]\tscate\tsgrp\n"; }
        else { print "\nError\n\n"; exit; }
    }
    print "\n\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;
}
}
} else {
    $prim =~ /\$spl[19]$/;
    if (($scate =~ /\completes/){($scate =~ /\INcom.(.*)$/)}{
        if ($grp =~ /\G1$/){
            se22g1+--; print RE "$num\t$$$pl[19]\tscate\tsgrp\n"; print LE22 "$num\n"; }
        elseif ($grp =~ /\G2$/){
            se22g2+--; print RE "$num\t$$$pl[19]\tscate\tsgrp\n"; print LE22 "$num\n"; }
        elseif ($grp =~ /\G3$/){
            se22g3+--; print RE "$num\t$$$pl[19]\tscate\tsgrp\n"; print LE22 "$num\n"; }
        elseif ($grp =~ /\G4$/){
            se22g4+--; print RE "$num\t$$$pl[19]\tscate\tsgrp\n"; print LE22 "$num\n"; }
        elseif ($grp =~ /\out_ranged.G(0-9)$/){
            se22out++; print RE "$num\t$$$pl[19]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\no.groups/) {
            se22out++; print RE "$num\t$$$pl[19]\tscate\tsgrp\n"; }
        else { print "\nError\n\n"; exit; }
    } elseif (($scate =~ /\OUT(.*)$/){
        if ($grp =~ /\G1$/){($oe22g1++); print RE "$num\tINSig_$$_$pl[19]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\G2$/){($oe22g2++); print RE "$num\tINSig_$$_$pl[19]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\G3$/){($oe22g3++); print RE "$num\tINSig_$$_$pl[19]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\G4$/){($oe22g4++); print RE "$num\tINSig_$$_$pl[19]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\out_ranged.G(0-9)$/){
            se22out++; print RE "$num\tINSig_$$_$pl[19]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\no.groups/) {
            se22out++; print RE "$num\tINSig_$$_$pl[19]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\noHPVs/) {($oe22nh++); print RE "$num\tINSig_$$_$pl[19]\tscate\tsgrp\n"; }
        else { print "\nError\n\n"; exit; }
    }
    print "\n\nThere is error in PrimerMatchingCategory...$num\n\n"; exit;
}
}
} else {
    $prim =~ /\$spl[20]$/;
    if (($scate =~ /\completes/){($scate =~ /\INcom.(.*)$/)}{
        if ($grp =~ /\G1$/){
            se23g1+--; print RE "$num\t$$$pl[20]\tscate\tsgrp\n"; print LE23 "$num\n"; }
        elseif ($grp =~ /\G2$/){
            se23g2+--; print RE "$num\t$$$pl[20]\tscate\tsgrp\n"; print LE23 "$num\n"; }
        elseif ($grp =~ /\G3$/){
            se23g3+--; print RE "$num\t$$$pl[20]\tscate\tsgrp\n"; print LE23 "$num\n"; }
        elseif ($grp =~ /\G4$/){
            se23g4+--; print RE "$num\t$$$pl[20]\tscate\tsgrp\n"; print LE23 "$num\n"; }
        elseif ($grp =~ /\out_ranged.G(0-9)$/){
            se23out++; print RE "$num\t$$$pl[20]\tscate\tsgrp\n"; }
        elseif ($grp =~ /\no.groups/) {
            se23out++; print RE "$num\t$$$pl[20]\tscate\tsgrp\n"; }
    }

```


File name: **set2_p6mac_fasta.pl**

Source code:

```
#!/usr/bin/perl

## About the script #####
# Created by Sasithorn Chotewutmontri, Jan 2009.
##
#####

use strict;
use warnings;

my $sampleNo = @ARGV;

my $prefixPath= $sampleNo[2];
my $pyroNum = int($sampleNo[1]);

my $path = $prefixPath."Pyro".$pyroNum."_result/";
my $fastaRC = $workingDir."sample".$sampleNo[0]."_RC.txt";

my $cut00 = "noCutoff";
my $cut28 = "28bpCutoff";

my @primerListEpyro3 = ("E17","E02","E03","E04","E05","E06","E07","E08",
"E09","E21","E11","E12","E13","E14","E15","E16",
"E18","E19","E20","E22","E23","E24","E25","E26",
"E27","E28","E29","E30","E31","E32");

print "\n\nAdding Fasta sequences after Program6....starts.....";

getFasta ($workingDir, $fastaRC, $cut00, \@primerListEpyro3);
getFasta ($workingDir, $fastaRC, $cut28, \@primerListEpyro3);

print "FINISHED\n\n";

exit;

##x Main Program ENDS HERE #####
#####

## SUBROUTINES #####
sub getFasta {
    my ($filename) = @_;
    use strict;
    use warnings;

    my @filedata = ();
    unless (open(FILEDATA, $filename)) {
        print STDERR "\n\nThe program can't open the files \"$filename\".\n\n";
        print "Please check input file name and its location\n";
        print "The correct command should be : \n\n";
        print "\t perl PROGRAM (locationname) \"$filename\"";
        print "\t perl PROGRAM (locationname) \"$filename\"";
        print "The command should be given under directory";
        print "C:/Perl/ bin> in case of Dos Terminal\n\n";
        exit;
    }
    @filedata = <FILEDATA>;
    close FILEDATA;
    return @filedata;
}

sub getFasta {
    my ($workDir, $fasInFile, $prefix, $prim) = @_;
    use strict;
    use warnings;

    my @inseq = getFasta ($fasInFile);
    my $dir = $workDir."intersection/".$prefix."_intersect/";
    my $surin = "_n_sighPV_list";
    my $surout = "_n_sighPV_fasta";
```

```
for (my $p = 0; $p < (scalar @prim); $p++) {
    my $listloc = $dir.$prefix.$prim[$p].$surin;
    my $outloc = $dir.$prefix.$prim[$p].$surout;
    my @list = getFasta ($listloc);
    unless (open (FA, ">$outloc")) {print "\n\nCan not create \"$outloc\".\n\n"; exit;}
    for (my $r = 0; $r < (scalar @list); $r++) {
        my $seqnum = int($list[$r]);
        my $nameLine = 2*$seqnum;
        my $seqLine = (2*$seqnum)+1;
        print FA $inseq[$nameLine].$inseq[$seqLine];
    }
    close FA;
}

}

##x Sub-routines END HERE #####
#####

File name: set2_p7mac.pl
Source code:

#!/usr/bin/perl

## About the script #####
# Created by Sasithorn Chotewutmontri, Jan 2009.
##
#####

## MAIN PROGRAM BODY #####
#####

use strict;
use warnings;

my @sampleNo = @ARGV;
my $prefixPath= $sampleNo[2];
my $pyroNum = int($sampleNo[1]);

my $path = $prefixPath."Pyro".$pyroNum."_result/";
my $samplePath = $path.$sampleNo[0]."/";
my $barcode = $sampleNo[0];

print "\n\nPROGRAM 7 (Find common 'breakpoint', if exists) STARTS.....\n\n";

my $cut00 = "noCutoff";
my $cut28 = "28bpCutoff";
my $cfcrcp_suffix = "_allreport_fullpieceposition";

my @primerListEpy3ro = ("E17","E02","E03","E04","E05","E06","E07","E08",
"E09","E21","E11","E12","E13","E14","E15","E16",
"E18","E19","E20","E22","E23","E24","E25","E26",
"E27","E28","E29","E30","E31","E32");

getBreakPoint2 ($samplePath, $cfcrcp_suffix, $cut00, \@primerListEpy3ro, $barcode);
getBreakPoint2 ($samplePath, $cfcrcp_suffix, $cut28, \@primerListEpy3ro, $barcode);

my $oldDirName = $samplePath."intersection/";
my $newDirName = $samplePath."intersection-".$pyroNum."_B".$sampleNo[0]."/";

system ("mv $oldDirName $newDirName");

print "FINISHED\n\n";
exit;

##x Main Program ENDS HERE #####
#####

## SUBROUTINES #####
#####

sub getFasta {
    my ($filename) = @_;
    use strict;
    use warnings;
```

```
my $score2 = $thit2[1];  
my $firstM2 = $thit2[2];  
my $strand2 = $thit2[3];  
my $firstMS2 = $thit2[4];  
  
$infoName = $info[2*$s];  
$infoSeq = ~ s/\$/g;  
$infoSeq = $info[(2*$s)+1];  
  
if ($seqnum == $num) {  
    my $pos0correct = $pos0 - 1 ;  
  
    print OUT $num, "\t", $tol, "\t", $pos0correct, "\t";  
    print ALL2 $bc, "\t", $s, $prjmlst[$p];  
    print ALL2 "\t", $num, "\t", $tol, "\t", $pos0correct, "\t";  
  
    if ($series1 =~ /^useAlones/) {  
        $posS2 = "none";  
        $posO2 = "none";  
        $BP2 = "none";  
        $overlap = "none";  
  
        if ($strand1 =~ /^plus/) {  
            $posS = int($firstMS1) + int($score1) - 1 ;  
        } elsif ($strand1 =~ /^minus/) {  
            $posS = int($firstMS1) - int($score1) + 1 ;  
        } else {  
            print "\n\nError occurs at strand name\n\n"; exit;  
        }  
    } elsif (($series1 =~ /^subsets/) {  
        $posS2 = "none";  
        $posO2 = "none";  
        $BP2 = "none";  
        $overlap = "none";  
  
        if ($strand2 =~ /^plus/) {  
            $posS = int($firstMS2) + int($score2) - 1 ;  
        } elsif ($strand2 =~ /^minus/) {  
            $posS = int($firstMS2) - int($score2) + 1 ;  
        } else {  
            print "\n\nError occurs at strand name\n\n"; exit;  
        }  
    }  
  
    } elsif (($series1 =~ /^use_1sts/) && ($series2 =~ /^use_2nds/)) {  
        $posO2 = "combined";  
        $BP2 = $score2;  
        my $difference = int($firstM1) + int($score1) - int($firstM2);  
        $overlap = "overlap=" . $difference;  
  
        if ($strand1 =~ /^plus/) {  
            $posS = int($firstMS1) + int($score1) - 1 ;  
        } elsif ($strand1 =~ /^minus/) {  
            $posS = int($firstMS1) - int($score1) + 1 ;  
        } else {  
            print "\n\nError occurs at strand name\n\n"; exit;  
        }  
    }  
  
    if ($strand2 =~ /^plus/) {  
        $posS2 = int($firstMS2) + int($score2) - 1 ;  
    } elsif ($strand2 =~ /^minus/) {  
        $posS2 = int($firstMS2) - int($score2) + 1 ;  
    } else {  
        print "\n\nError occurs at strand name\n\n"; exit;  
    }  
}  
  
} elsif (($series1 =~ /^use_1sts/) && ($series2 =~ /^use_2nd_discounts/)) {  
    $posO2 = int($firstM2) + int($score2) - 1;  
    $BP2 = $score2;  
    $overlap = "none";  
  
    if ($strand1 =~ /^plus/) {  
        $posS = int($firstMS1) + int($score1) - 1 ;  
    } elsif ($strand1 =~ /^minus/) {  
        $posS = int($firstMS1) - int($score1) + 1 ;  
    } else {  
        print "\n\nError occurs at strand name\n\n"; exit;  
    }  
}
```

```

unless ( $? ) {
    my @fileData = (FILE_DATA, $filename) ) {
        print STDERR "The program can't open the files \"$filename\"\\n\\n\\n";
        print "Please re-check input file name and its location\\n";
        print "The correct command should be: \\n\\n";
        print "\\tperl PROGRAM(location-name) INPUTFASTFILE(location+name)\\n\\n";
        print "The command should be given under directory:\\n";
        print " c:/Perl/bin/ in case of Dos Terminal\\n\\n";
        exit;
    }
}
@fileData = <FILE_DATA>;
close FILE_DATA;
return @fileData;
}

sub getBreakPoint2 {
    my ($path, $repSuffix, $prefix, $primList, $bc) = @_;
    use strict;
    use warnings;

    my $subPathCf = $path.$prefix."/";
    my $CFRepFile = $subPathCf.$prefix."_allreport_fullpieceposition";
    my @incf = getFileData ($CFRepFile);

    my $subPathInter = $path."intersection/".$prefix."_intersect/";
    my $suffix = "_n_sighpv.list";
    my $sufas = "_n_sighpv.fasta";
    my $soutall2 = $subPathInter."PrimerAll_breakpoint_".$prefix."_2";

    unless (open (ALL2, ">$soutall2")) {print "\\n\\n\\nCan not create \"$soutall2\"\\n\\n\\n"; exit;}}

    print ALL2 "In ".$prefix."\\n\\n";
    print ALL2 "Barcode\\tPrimer\\tSeqNo\\tSeqLength\\tLastPositionQuery\\tLastPositionHPV\\n";
    print ALL2 "\\tMatchBP\\n";
    print ALL2 "\\tStartPositionQuery_2ndHit\\tLastPositionHPV_2ndHit\\n";
    print ALL2 "\\tMatchBP_2ndHit\\tOverlap\\tOriginaName\\tSequence\\n";

    # for each primer
    for (my $p = 0; $p < (scalar @primList); $p++) {

        my $inListFile = $subPathInter.$prefix.$primList[$p].$suffix;
        my $infasFile = $subPathInter.$prefix.$primList[$p].$sufas;
        my $sout = $subPathInter.$primList[$p]."_breakpoint_".$prefix;

        unless (open (OUT, ">$sout")) {print "\\n\\n\\nCan not create \"$sout\"\\n\\n\\n"; exit;}}

        print OUT "For-PRIMER GROUP ".$primList[$p]. "\\n";
        print OUT "\\tSeqNo\\tSeqLength\\tLastPositionQuery\\tLastPositionHPV\\n";
        print OUT "\\tMatchBP\\tOriginaName\\tSequence\\n";

        my @inList = getFileData ($inListFile);
        my @infas = getFileData ($infasFile);

        for (my $a = 0; $a < (scalar @inList); $a++) {
            my $seqnum = int($inList[$a]);

            my @temp_cfln = split (/#/, $inCf[$seqnum]);
            my $stnum = $temp_cfln[0];
            my $stOL = $temp_cfln[1];

            my @temp_block = split (/%, $temp_cfln[3]);
            my $tBP = $temp_block[0];
            my $tMS = $temp_block[3];
            my $tPos0 = $temp_block[4];
            my ($stpos5, $infasName, $infasSeq);
            my ($stpos02, $stpos52, $tBP2, $soverlap);

            my $usage = $temp_cfln[4];
            # indicates which hit(s) is used AND how it's used

            my @thit1 = split (/%, $temp_cfln[5]);
            my $series1 = $thit1[0];
            my $score1 = $thit1[1];
            my $firstM1 = $thit1[2];
            my $strand1 = $thit1[3];
            my $firstMS1 = $thit1[4];

            my @thit2 = split (/%, $temp_cfln[6]);
            my $series2 = $thit2[0];

```

File name: **set3_p8mac.pl**

Source code:

```
#!/usr/bin/perl

## About the script #####
# Created by Sasithorn Chotevitmontri, Jan 2009.
#
#####
### MAIN PROGRAM BODY #####
#####
use strict;
use warnings;

my @sampleNo = @ARGV;
my $prefixPath= $sampleNo[2];
my $pyRunNo = int($sampleNo[1]);

my $path = $prefixPath."Pyro".$pyRunNo."_result/";
my $clus_loc = "/clustalw-2.0.10-macosx/clustalw2"; # path of clustalw2
my $samplePath = $path.$sampleNo[0]."/";
my $barcode = $sampleNo[0];
my $barcodeB = $pyRunNo."B".$barcode;
my $hpvEE_file = ">". $pyRunNo."B".$barcode;
my $hpvEE_path = "PERL/ForPyro3/hpv16R_EEprimer_15Int"; # path of EEprimer_15Int
my $cut00 = "noCutoff";
my $cut28 = "28bpCutoff";

print "\n\nPROGRAM 8 (multiple alignment, automatic) STARTS.....";

my @hpvEE = getFileData ($hpvEE_file);
my @primerListEep3ro = ("E17", "E02", "E03", "E04", "E05", "E06", "E07", "E08",
"E09", "E21", "E11", "E12", "E13", "E14", "E15", "E16",
"E18", "E19", "E20", "E22", "E23", "E24", "E25", "E26",
"E27", "E28", "E29", "E30", "E31", "E32");

my @hpvPrimerEeposP3ro = ("810", "1275", "1576", "1951", "2403", "2723", "3121", "3555",
"1064", "1471", "1785", "2174", "2569", "2933", "3339", "3778",
"968", "1115", "1361", "1672", "1860", "2069", "2288", "2476",
"2628", "2857", "3037", "3199", "3455", "3696");

# (1) Prepare FASTA-formatted SEQUENCES OF EACH BARCODE. EACH PRIMER
getPrepSeq ($samplePath, $cut00, \@primerListEep3ro, $barcode, $fname1, \@hpvEE, \@hpvPrimerEeposP3ro);
getPrepSeq ($samplePath, $cut28, \@primerListEep3ro, $barcode, $fname1, \@hpvEE, \@hpvPrimerEeposP3ro);

# (2) ALIGN using CLUSTAL W, and get OUTPUTS (in FASTA format) into a new directory
getMultiAlignClustalW ($samplePath, $cut00, \@primerListEep3ro, $barcode, $clus_loc);
getMultiAlignClustalW ($samplePath, $cut28, \@primerListEep3ro, $barcode, $clus_loc);

# (3) MOVE the dnd OUTPUTS to the same directory as the ALIGNED sequence outputs
getDndMoved ($samplePath, $cut00, \@primerListEep3ro, $barcode);
getDndMoved ($samplePath, $cut28, \@primerListEep3ro, $barcode);

print "FINISHED\n";
exit;

#xx Main Program ENDS HERE
#####
##### SUBROUTINES #####
#####
sub getFileData {
    my ($filename) = @_;
    use strict;
    use warnings;

    my @fileData = ();
    unless ( open(FILE_DATA, $filename) ) {
        print STDERR "\n\nThis program can't open the files \"$filename\".\n\n";
        print "Please re-check input file name and its location\n";
        print "\n\nThe correct command should be : perl PROGRAM(location-name) INPUTFASTFILE(location+name)\n\n";
        print "The command should be given under directory;\n";
        print " c:/perl/bin/ in case of dos terminal\n\n";
    }
}

#####
##### Sub-routines END HERE #####
#####
```

```

    exit;
}
@filedata = <FILE_DATA>;
close FILE_DATA;
return @filedata;
}

sub getPrepSeq {
    my ($path, $prefix, $primList, $bc, $fname1, $hpvee, $hpvPosArray) = @_;
    use strict;
    use warnings;

    my $subPathInter = $path."intersection-".$bc."/". $prefix."_intersect/";

    # for each primer
    for (my $p = 0; $p < (scalar @primList); $p++) {
        my $inPrefix = $bc." ".$primList[$p];
        my $inPath = $subPathInter.$inPrefix."_editFasta_Hpvee_Selected_".$prefix;
        my $outPath = $subPathInter.$inPrefix."_editFasta_Hpvee_Selected_".$prefix."_ALIGNED";
        my $command = $clus_loc." -INFILE=".$inPath." -ALIGN -OUTFILE=".$outPath." -OUTPUT=FASTA";
        system ("command");
    }
}

sub getDndMoved {
    my ($path, $prefix, $primList, $bc) = @_;
    use strict;
    use warnings;

    my $subPathInter = $path."intersection-".$bc."/". $prefix."_intersect/";
    my $outDir = $path."ALIGNMENTS-".$bc."/";

    # for each primer
    for (my $p = 0; $p < (scalar @primList); $p++) {
        my $inPrefix = $bc." ".$primList[$p];
        my $inPath = $subPathInter.$inPrefix."_editFasta_Hpvee_Selected_".$prefix;
        my $command = "mv ".$inPath."_dnd ".$outDir;
        system ("command");
    }
}

##xx Sub-routines END HERE #####
#####

#####
File name:      set3_p9mac.p1
Source code:

#!/usr/bin/perl

## About the script #####
##
## Created by Sasithorn Chotewitmontri, Jan 2009.
##
#####
##### MAIN PROGRAM BODY #####
#####
#####
use strict;
use warnings;

my @sampleNo = @ARGV;
my $prefixPath = $sampleNo[2];
my $pyRunNo = int($sampleNo[1]);

my $path = $prefixPath."Pyro". $pyRunNo."_result/";
my $samplePath = $path.$sampleNo[0]."/";
my $barcode = $sampleNo[0];

my $barcodeB = $pyRunNo."B". $barcode;
my $cut00 = "nocutoff";
my $cut28 = "28bcutoff";

print "\n\nPROGRAM 9 (sort ALIGNED-fasta sequences acc to size) STARTS.....";

my @primerListEpy3ro = ("E17", "E02", "E03", "E04", "E05", "E06", "E07", "E08",
                        "E09", "E11", "E12", "E13", "E14", "E15", "E16",
                        "E18", "E19", "E20", "E21", "E22", "E23", "E24", "E25", "E26",
                        "E27", "E28", "E29", "E30", "E31", "E32");

```

```

getALignedSeqSorted ($samplePath, $cut28, \@primerListEpy3ro, $barcodeB);
getALignedSeqSorted ($samplePath, $cut00, \@primerListEpy3ro, $barcodeB);

print "FINISHED\n\n";
exit;

##xx Main Program ENDS HERE xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
#####
### SUBROUTINES #####
#####
sub getFastData {
    my ($filename) = @_;
    use strict;
    use warnings;

    my @filedata = ();
    unless ( open(FILE_DATA, $filename) ) {
        print STDERR "\n\nThe program can't open the files \"$filename\"(\"$location\n";
        print "Please re-check input file name and its location\n";
        print "The correct command should be :\n\n";
        print "\t\tperl PROGRAM(location+name) INPUTFASTFILE(location+name)\n\n";
        print "The command should be given under directory";
        print " : c:/Perl/bin> in case of Dos Terminal\n\n\n";
        exit;
    }
    @filedata = <FILE_DATA>;
    close FILE_DATA;
    return @filedata;
}

sub getALignedSeqSorted {
    my ($path, $prefix, $primList, $bc) = @_;
    use strict;
    use warnings;

    my $indir = $path."ALIGNMENTS-".$bc."/";

    # for each primer
    for (my $p = 0; $p < (scalar @primList); $p++) {

        my $inprefix = $bc."_"."$primList[$p];

        my $input = $indir.$inprefix."_editFA.HVVEE_selected_".$prefix."_ALIGNED";
        my $output = $indir.$inprefix."_editFA.HVVEE_selected_".$prefix."_ALIGNED_unsorted";
        my $output2 = $indir.$inprefix."_editFA.HVVEE_selected_".$prefix."_ALIGNED_SORTED";
        my $output3 = $indir.$inprefix."_editFA.HVVEE_selected_".$prefix."_ALIGNED_SORTED_Fasta";

        my @infas = getFastData ($input);
        my @STACK2 = extractFastStack (@infas);
        my @fname = ();
        my @fseq = ();
        extractFastName2 (\@infas, \@STACK2, \@fname);
        extractFastSeq2 (\@infas, \@STACK2, \@fseq);

        unless (open (OUT, ">$output")) {print "\n\n\nCan not create \"$output\".\n\n\n"; exit;}
        unless (open (OUT2, ">$output2")) {print "\n\n\nCan not create \"$output2\".\n\n\n"; exit;}
        unless (open (OUT3, ">$output3")) {print "\n\n\nCan not create \"$output3\".\n\n\n"; exit;}

        print OUT "name\tsequence\tgap-counted\tindex\n";
        print OUT2 "name\tsequence\tgap-counted\tindex\n";

        my $index;
        my $seq;
        my $name;
        my @sortedN_keys;
        # define array of the SeqName after sorting by 'gapcount'

        # for each fasta seq
        for (my $m = 0; $m < (scalar @fname); $m++) {

            my ($gapstring, $gapcount);

            my $modiName = $fname[$m];
            my $modiSeq = $fseq[$m];

            $modiName =~ s/\s//g;
            $modiSeq =~ s/\s//g;

            # fasta name, still with a new line at the end
            # fasta seq, still with new-lines

            # get rid of 'new line' character(s)
            # get rid of 'new line' character(s)

```

```

        $N_index($modiName) = "$m"; # input the 'value' (index) for 'key' (SeqName) for hash %N_index
        $N_seq($modiName) = "$modiSeq"; # input the 'value' (Seq) for 'key' (SeqName) for hash %N_seq

        # control the 'ending' part of the Seq -- should be '-' character
        # IF it's the longest seq and has NO '-' character at the end,
        # the 'gapcount' for this case was given as '0' by the program
        if ( $modiSeq =~ /[A-Z]((-)*$)/1 ) {

            $gapstring = $1;
            $gapcount = length ($gapstring); # count the number of '-' character
            print OUT $modiName."\\t".$modiSeq."\\t".$gapcount."\\t".$m."\\n";

            # input the 'value' (gapcount) for 'key' (SeqName) for hash %N_gapcount
            $N_gapcount($modiName) = "$gapcount";

        }

        my @keysN = keys %N_gapcount; # unsorted keys of the hash %N_gapcount, put into an array @keysN

        # sorted keys of the hash %N_gapcount (sorted by its VALUES, descending) are put in array @sortedN_keys
        # i.e. sort function { using condition: hash-values (descending) } of the keys of the hash
        # and then return the keys of the hash
        # --> $a and $b represent in this case the keys
        @sortedN_keys = sort { $N_gapcount{$a} <=> $N_gapcount{$b} } keys %N_gapcount;

        # for each value of the array @sortedN_keys which contains the sorted SeqNames
        # the keys (SeqName) are called to get the corresponding values (Seq) from hash %N_seq
        for (my $n = 0; $n < (scalar @sortedN_keys); $n++) {
            print OUT2 $sortedN_keys[$n]."\\t".$N_seq{$sortedN_keys[$n]}."\\n";
            print OUT3 $sortedN_keys[$n]."\\n".$N_seq{$sortedN_keys[$n]}."\\n";
        }

        close OUT; close OUT2; close OUT3;

    }

}

sub extractFastStack {
    my (@fastFileData) = @_;
    use strict;
    use warnings;

    my $nmaxcount = scalar @fastFileData;
    my @stacklist = ();
    for (my $m = 0; $m < $nmaxcount; $m++) {
        if ($fastFileData[$m] =~ /\s*$/) {
            push @stacklist, 2;
        } elsif ($fastFileData[$m] =~ /\s*$/) {
            push @stacklist, 2;
        } elsif ($fastFileData[$m] =~ /\s*$/) {
            push @stacklist, 0;
        } else {
            push @stacklist, 1;
        }
    }
    return @stacklist;
}

sub extractFastName2 {
    my ($data1, $STACK, $data1_name) = @_;
    use strict;
    use warnings;

    my $nmaxcount = scalar @$data1;
    for (my $p = 0; $p < $nmaxcount; $p++) {
        if ($STACK[$p] =~ /\d/) {
            push @$data1_name, $data1[$p];
        } else {
            next;
        }
    }
}

sub extractFastSeq2 {
    my ($data1, $STACK, $data1_seq) = @_;
    use strict;
    use warnings;

```

```
#####
my $maxcounter = scalar @data1;
for (my $p=0; $p < $maxcounter ; $p++) {
    if ($$STACK[$p] =~ "0") {
        } else if ( ($$STACK[$p] =~ "1") && ($$STACK[$p - 1] =~ "0")) {
            push @data1_seq, $$data1[$p];
        } elsif (($$STACK[$p] =~ "1") && ($$STACK[$p - 1] =~ "1")) {
            my $lastEntry = (scalar @data1_seq) - 1;
            $$data1_seq[$lastEntry] .= $$data1[$p];
        } else { next; }
    }
}

#####
####xx Sub-routines END HERE #####
#####

File name:          set4_p10mac.p1
Source code:

#!/usr/bin/perl

#####
## About the script #####
#####
# Created by Sasithorn Chotewutmontri, Jan 2009.
#####
#### MAIN PROGRAM BODY #####
#####
#####
#####
use strict;
use warnings;

my @sampleNo = @ARGV;

my $input = $sampleNo[0];
my $prefix_path = $sampleNo[2];
my $pyRunNo = int($sampleNo[1]);

print "\n\n===== \n\n";
print "Running PROGRAM 10, FOR ALL BARCODES, ASP14 No: ". $pyRunNo. "\n\n";

my $outpath = $prefix_path."Pyro".$pyRunNo."_result_ALIGN/";

my @primerListEepY3ro = ("E17","E02","E03","E04","E05","E06","E07","E08",
"E09","E21","E11","E12","E13","E14","E15","E16",
"E18","E19","E20","E22","E23","E24","E25","E26",
"E27","E28","E29","E30","E31","E32");

system ("mkdir", $outpath);

getFolders (@primerListEepY3ro, $outpath);

print "FINISHED\n\n";
exit;

####xx Main Program BNDs HERE #####
#####
##### Sub-routine #####
#####
sub getFolders {
    my ($primList, $outpathprefix) = @_;
    use strict;
    use warnings;
    # for each primer
    for (my $p=0; $p < (scalar @$primList); $p++) {
        my $newDir = $outpathprefix."/".$$primList[$p]."/";
        system ("mkdir", $newDir);
    }
}

####xx Sub-routines END HERE #####
#####
```